# 07 - Bootstrap and Splines

Data Mining

*SYS 6018 | Fall 2019*

*07-bootstrap.pdf*

## Contents

# 1   Introduction to the Bootstrap

## 1.1   Required R Packages

We will be using the R packages of:

- `boot` for help with bootstrapping
- `broom` for tidy extraction of model components
- `splines` for working with B-splines
- `tidyverse` for data manipulation and visualization

```r
library(boot)
library(broom)
library(splines)
library(tidyverse)
```

## 1.2   Uncertainty in a test statistic

There is often interest to understand the uncertainty in the estimated value of a test statistic.

- For example, let $p$ be the actual/true proportion of customers who will use your company's coupon.

- To estimate $p$, you decide to take a sample of $n = 200$ customers and find that $x = 10$ or $\hat{p} = 10/200 = 0.05 = 5\%$ redeemed the coupon.
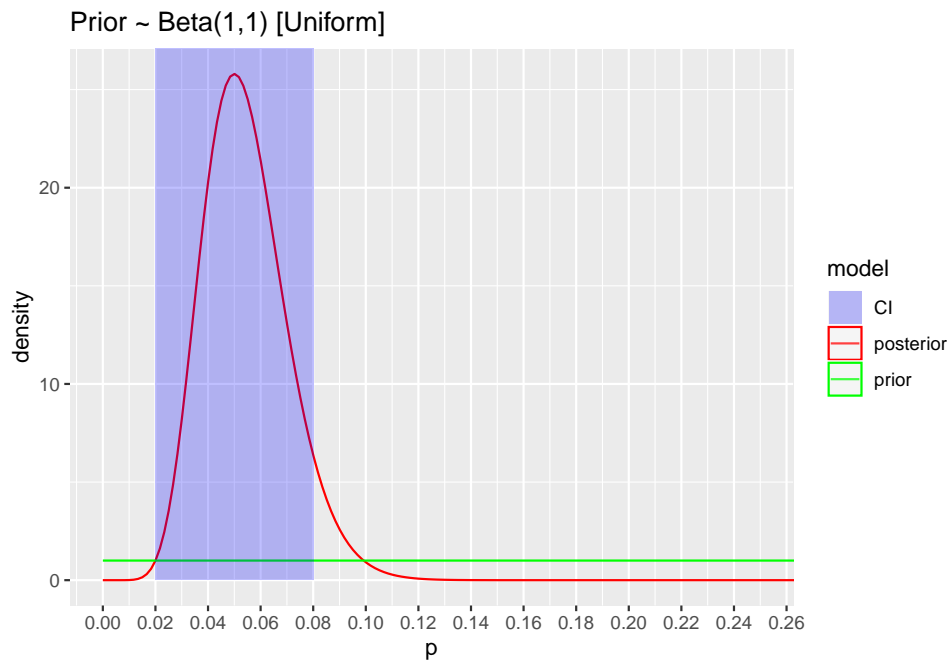
### 1.2.1   Confidence Interval

- It is common to calculate the 95% *confidence interval (CI)*

$$\mathrm{CI}(p) = \hat{p} \pm 2 \cdot \mathrm{SE}(\hat{p})$$

$$= \hat{p} \pm 2\sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

$$= 0.05 \pm 0.03$$

- This calculation is based on the assumption that $\hat{p}$ is approximately normally distributed with the mean equal to the *unknown* true $p$, i.e., $\hat{p} \sim N(p, \sqrt{\frac{p(1-p)}{n}})$.
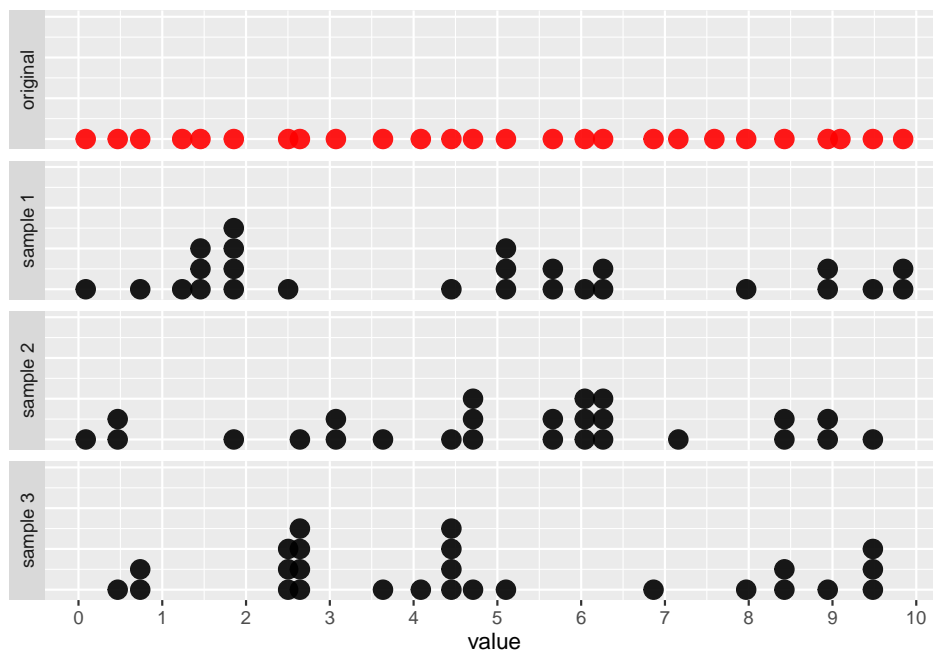
### 1.2.2   Bayesian Posterior Distribution

In the Bayesian world, you'd probably specify a Beta *prior* for $p$, i.e., $p \sim \mathrm{Beta}(a, b)$ and calculate the *posterior* distribution $p \mid x = 10 \sim \mathrm{Beta}(a + x, b + n - x)$ which would fully characterize the uncertainty.

Prior ~ Beta(1,1) [Uniform]

### 1.2.3 The Bootstrap

- The Boostrap is a way to assess the uncertainty in a test statistic using *resampling*.

- The idea is to simulate the data from the *empirical distribution*, which puts a point mass of $1/n$ at each observed data point (i.e., sample the original data **with replacement**).

    - It is important to simulate $n$ observations (same size as original data) because the uncertainty in the test statistic is a function of $n$

- Then, calculate the test statistic for each bootstrap sample. The variability in the collection of bootstrap test statistics should be similar to the variability in the test statistic.

---

### Algorithm: Nonparametric/Empirical Bootstrap

Observe data $D = [X_1, X_2, \ldots, X_n]$ ($n$ observations).
Calculate a test statistic $\hat{\theta} = \hat{\theta}(D)$, which is a function of $D$.
Repeat steps 1 and 2 $M$ times:
  1. Simulate $D^*$, a new data set of $n$ observations by sampling from $D$ with replacement.
  2. Calculate the bootstrap test statistic $\hat{\theta}^* = \hat{\theta}(D^*)$
The bootstrapped samples $\hat{\theta}_1^*, \hat{\theta}_2^*, \ldots, \hat{\theta}_M^*$ can be used to estimate the distribution of $\hat{\theta}$.
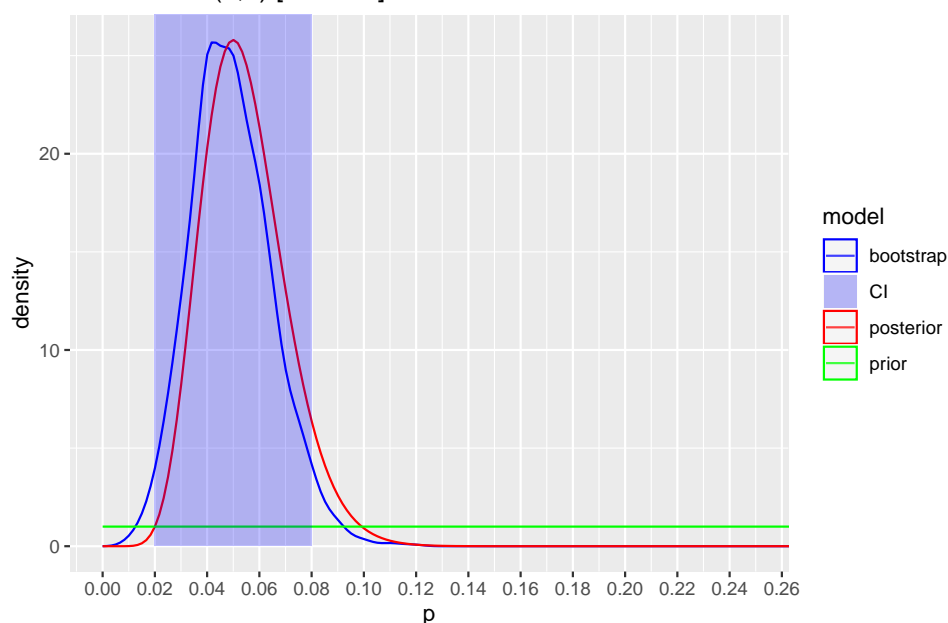  - Or properties of the distribution, like standard deviation (standard error), percentiles, etc.

---

```r
#-- Original Data
x = c(rep(1, 10), rep(0, 190))      # 10 successes, 190 failures
n = length(x)                        # length of observed data

#-- Bootstrap Distribution
M = 2000                             # number of bootstrap samples
p = numeric(M)                       # initialize vector for test statistic
set.seed(201910)                     # set random seed
for(m in 1:M){
  #- sample from empirical distribution
  ind = sample(n, replace=TRUE)      # sample indices with replacement
  xboot = x[ind]                     # bootstrap sample
  #- calculate proportion of successes
  p[m] = mean(xboot)                 # calculate test statistic
}

#-- Bootstrap Percentile based confidence Intervals
quantile(p, probs=c(.025, .975))    # 95% bootstrap interval
#>  2.5% 97.5%
#>  0.02  0.08
```
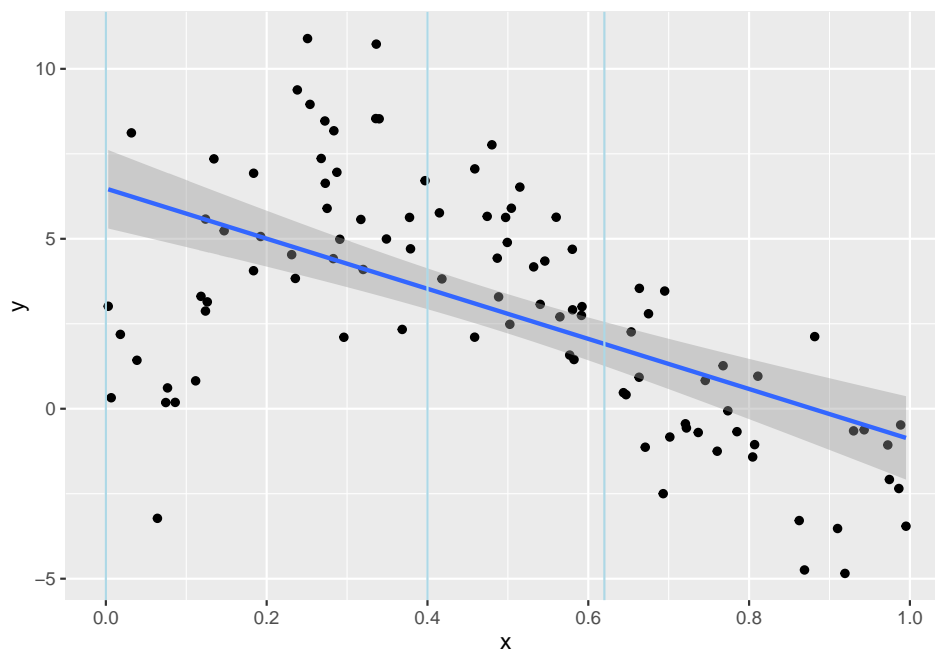


Prior ~ Beta(1,1) [Uniform]

- Notice that in the above example the bootstrap distribution is close to the Bayesian posterior distribution (using the uninformed Uniform prior).

- This is no accident, it turns out there is a close correspondence between the bootstrap derived distribution and the Bayesian posterior distribution under *uninformative priors*
    - See ESL 8.4 for more details

## 2    Bootstrapping Regression Parameters

The bootstrap is not limited to univariate test statistics. It can be used on multivariate test statistics.

Consider the uncertainty in estimates of the parameters (i.e., $\beta$ coefficients) of a regression model.



```
1  data_train = tibble(x,y)        # create a data frame/tibble
2  m1 = lm(y~x, data=data_train)  # fit simple OLS
3  broom::tidy(m1, conf.int=TRUE)  # OLS estimated coefficients
4  #> # A tibble: 2 x 7
5  #>   term        estimate std.error statistic  p.value conf.low conf.high
6  #>   <chr>          <dbl>     <dbl>     <dbl>    <dbl>    <dbl>     <dbl>
7  #> 1 (Intercept)     6.48     0.584      11.1  5.39e-19     5.32      7.64
8  #> 2 x              -7.37      1.06     -6.97 3.69e-10    -9.47     -5.27
```

```
1  vcov(m1)                        # variance matrix
2  #>             (Intercept)         x
3  #> (Intercept)      0.3414  -0.5359
4  #> x               -0.5359   1.1183
```

### 2.1    Bootstrap the $\beta$'s
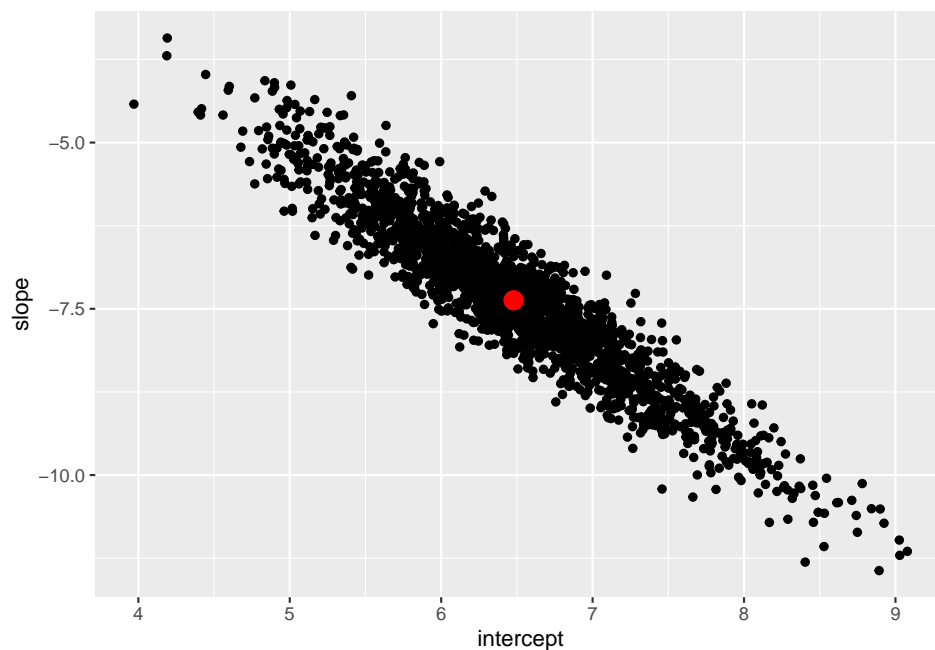
```r
1   #-- Bootstrap Distribution
2   M = 2000                            # number of bootstrap samples
3   beta = matrix(NA, M, 2)             # initialize vector for test statistics
4   set.seed(201910)                    # set random seed
5   for(m in 1:M){
6     #- sample from empirical distribution
7     ind = sample(n, replace=TRUE)     # sample indices with replacement
8     data.boot = data_train[ind,]      # bootstrap sample
9     #- fit regression model
10    m.boot = lm(y~x, data=data.boot) # fit simple OLS
11    #- save test statistics
12    beta[m, ] = coef(m.boot)
13  }
14  #- convert to tibble (and add column names)
15  beta = as_tibble(beta, .name_repair = "unique") %>%
16    setNames(c('intercept', 'slope'))
17
18  #- Plot
19  ggplot(beta, aes(intercept, slope)) + geom_point() +
20    geom_point(data=tibble(intercept=coef(m1)[1],
21                           slope = coef(m1)[2]), color="red", size=4)
22
23  #- bootstrap estimate
24  var(beta)                # varaince matrix
25  #>            intercept    slope
26  #> intercept     0.5821 -0.8747
27  #> slope        -0.8747  1.4918
```

```r
1   apply(beta, 2, sd)    # standard errors (sqrt of diagonal)
2   #> intercept      slope
3   #>    0.7629     1.2214
```

# 3 Basis Function Modeling

For a univariate $x$, a linear basis expansion is

$$\hat{f}(x) = \sum_j \hat{\theta}_j b_j(x)$$

where $b_j(x)$ is the value of the $j$th basis function at $x$ and $\theta_j$ is the coefficient to be estimated.

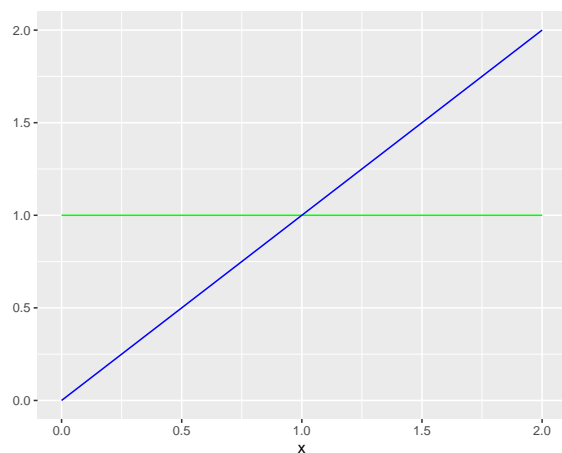- The $b_j(x)$ are usually specified in advanced (i.e., not estimated)

Examples:

- **Linear Regression**

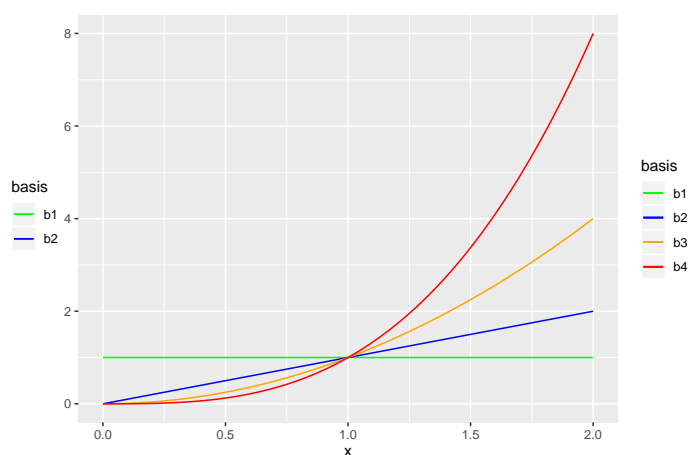$$\hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$b_1(x) = 1$$
$$b_2(x) = x$$

- **Polynomial Regression**

$$\hat{f}(x) = \sum_{j=1}^{d} \hat{\beta}_j x^j$$

$$b_j(x) = x^j$$

## 3.1   Piecewise Polynomials

Piecewise Constant

Piecewise Linear

Continuous Piecewise Linear

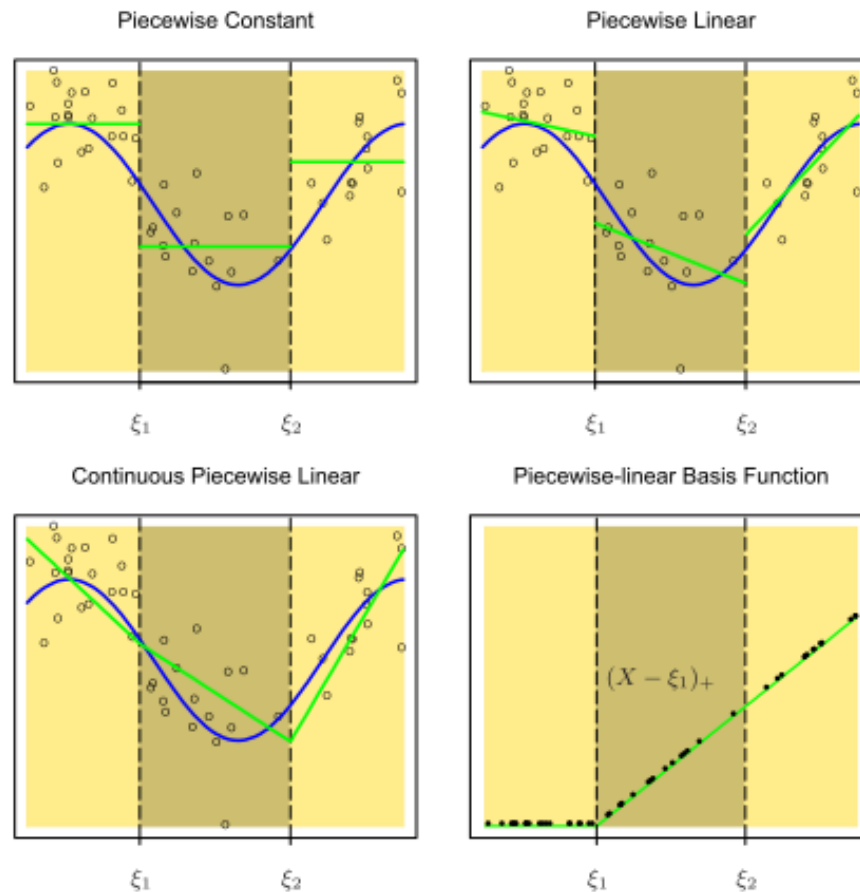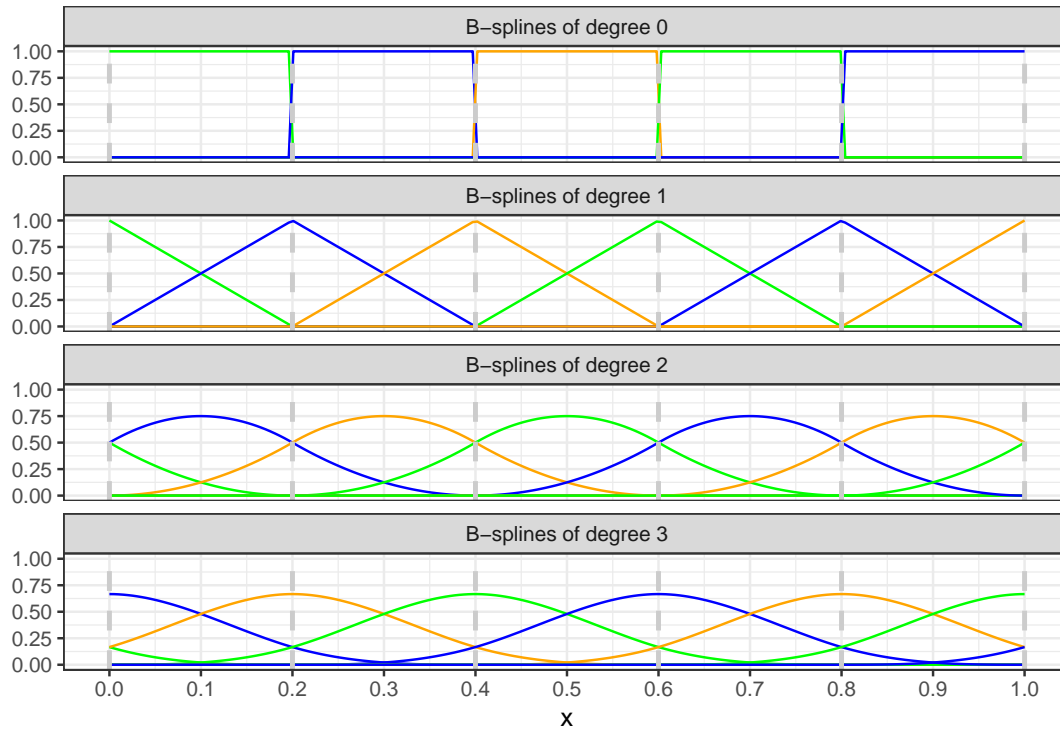Piecewise-linear Basis Function

$(X - \xi_1)_+$

**FIGURE 5.1.** *The top left panel shows a piecewise constant function fit to some artificial data. The broken vertical lines indicate the positions of the two knots $\xi_1$ and $\xi_2$. The blue curve represents the true function, from which the data were generated with Gaussian noise. The remaining two panels show piecewise linear functions fit to the same data—the top right unrestricted, and the lower left restricted to be continuous at the knots. The lower right panel shows a piecewise–linear basis function, $h_3(X) = (X - \xi_1)_+$, continuous at $\xi_1$. The black points indicate the sample evaluations $h_3(x_i)$, $i = 1, \ldots, N$.*

## 3.2   B-Splines



Like ESL Fig 5.20, B-splines (knots shown by vertical dashed lines)

- A degree = 0 B-spline is a *regressogram* basis. Will lead to a piecewise constant fit.

- A degree = 3 B-spline (called *cubic* splines) is similar in shape to a Gaussian pdf. But the B-spline has finite support and facilitates quick computation (due to the induced sparseness).

### 3.2.1   Parameter Estimation

In matrix notation,

$$\hat{f}(x) = \sum_j \hat{\theta}_j b_j(x)$$

$$= B\hat{\theta}$$

where $B$ is the *basis matrix*.

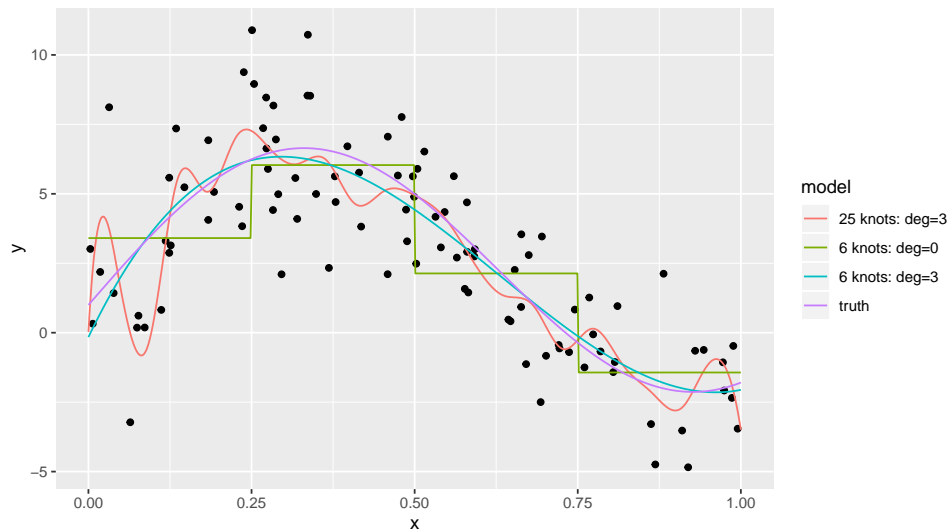- For example, a polynomial matrix is

$$B = \begin{bmatrix} 1 & X_1 & \ldots & X_1^J \\ 1 & X_2 & \ldots & X_2^J \\ \vdots & \vdots & \vdots & \vdots \\ 1 & X_n & \ldots & X_n^J \end{bmatrix}$$

- More generally,

$$B = \begin{bmatrix} b_1(x_1) & b_2(x_1) & \dots & b_J(x_1) \\ b_1(x_2) & b_2(x_2) & \dots & b_J(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ b_1(x_n) & b_2(x_n) & \dots & b_J(x_n) \end{bmatrix}$$

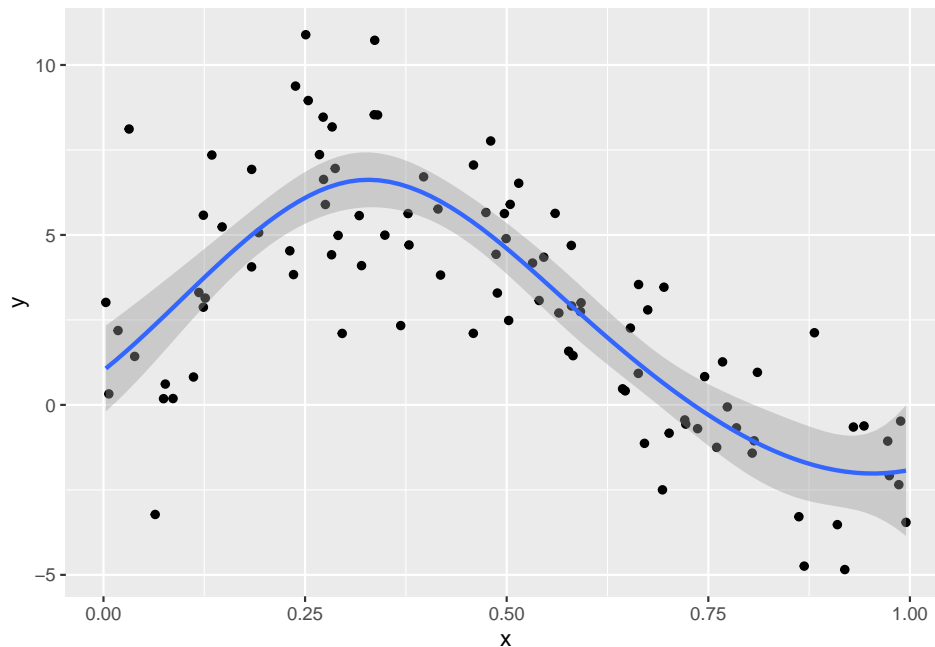- Now, its in a form just like linear regression! Estimate with OLS

$$\hat{\theta} = (B^\mathsf{T} B)^{-1} B^\mathsf{T} Y$$



- It may be helpful to think of a basis expansion as similar to a dummy coding for categorical variables.
    - This expands the single variable $x$ into $df$ new variables.
- In R, the function bs() can be put directly in formula to make a B-spline.

```r
#- fit a 5 df B-spline
# Note: don't need to include an intercept in the lm()
# Note: the boundary.knots are set just a bit outside the range of the data
#        so prediction is possible outside the range (see below for usage),
#        and will lead to equal knot locations during bootstrap.
#        You probably won't need to set this in practice.
kts.bdry = c(-.2, 1.2)
model_bs = lm(y~bs(x, df=5, deg=3, Boundary.knots = kts.bdry)-1,
              data=data_train)
tidy(model_bs)
#> # A tibble: 5 x 5
#>   term                                   estimate std.error statistic  p.value
#>   <chr>                                     <dbl>     <dbl>     <dbl>    <dbl>
#> 1 bs(x, df = 5, deg = 3, Boundary.knots =~  -2.50      1.51     -1.65  1.02e- 1
#> 2 bs(x, df = 5, deg = 3, Boundary.knots =~  10.9       1.27      8.61  1.53e-13
#> 3 bs(x, df = 5, deg = 3, Boundary.knots =~  -0.241     1.53     -0.157 8.76e- 1
#> 4 bs(x, df = 5, deg = 3, Boundary.knots =~  -4.71      3.07     -1.53  1.28e- 1
#> 5 bs(x, df = 5, deg = 3, Boundary.knots =~   1.45      6.90      0.211 8.34e- 1
```

```r
ggplot(data_train, aes(x,y)) + geom_point() +
  geom_smooth(method='lm', formula='y~bs(x, df=5, deg=3, Boundary.knots = kts.bdry)-1')
```
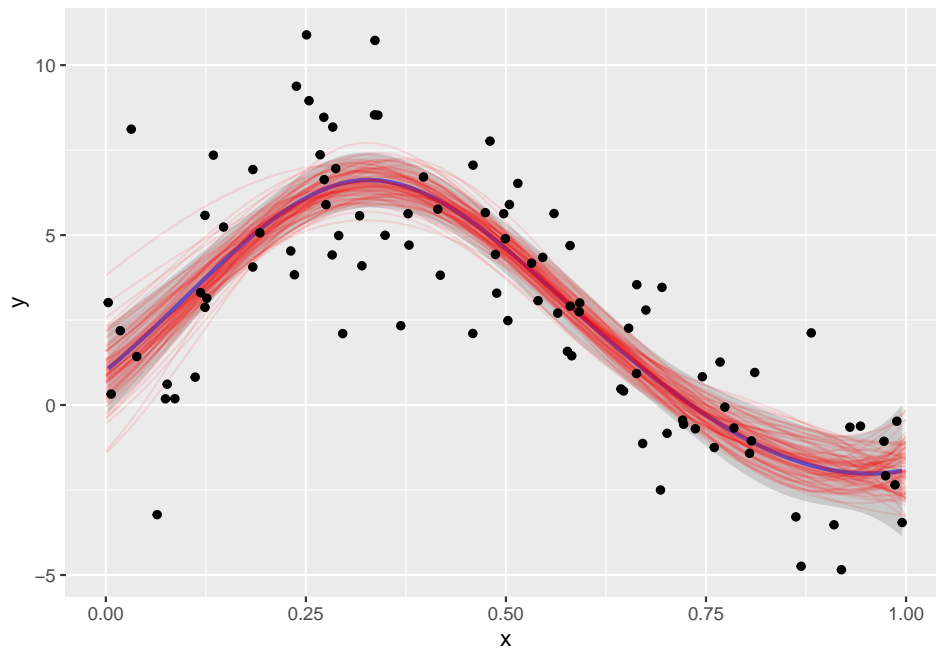
- Setting `df=5` will create a B-spline design matrix with 5 columns
  - So there are 5 basis functions
  - The equation will work with a non-integer df, but its not clear what it means. (Note: we will encounter *penalized* splines later when we can set a real-valued df, but `bs()` is not doing this.)
- The number of (internal) knots is equal to df-degree and at equally spaced quantiles of the data
  - With `df=5` and `deg=3`, there are 2 internal knots at the 33.33% and 66.66% percentiles of $x$

## 3.3   Bootstrap Confidence Interval for $f(x)$

Bootstrap can be used to understand the uncertainty in the fitted values

```
1   #-- Bootstrap CI (Percentile Method)
2   M = 100                                      # number of bootstrap samples
3   data_eval = tibble(x=seq(0, 1, length=300))  # evaluation points
4   YHAT = matrix(NA, nrow(data_eval), M)         # initialize matrix for fitted values
5
6   #-- Spline Settings
7   for(m in 1:M){
8     #- sample from empirical distribution
9     ind = sample(n, replace=TRUE)                  # sample indices with replacement
10    #- fit bspline model
11    m_boot = lm(y~bs(x, df=5, Boundary.knots=kts.bdry)-1,
12               data=data_train[ind,])   # fit bootstrap data
13    #- predict from bootstrap model
14    YHAT[,m] = predict(m_boot, newdata=data_eval)
15  }
16
17  #-- Convert to tibble and plot
18  data_fit = as_tibble(YHAT) %>%
19    bind_cols(data_eval) %>%     # add the eval points
20    gather(simulation, y, -x)     # convert to long format
21
22  ggplot(data_train, aes(x,y)) +
```

```
23    geom_smooth(method='lm',
24               formula='y~bs(x, df=5, deg=3, Boundary.knots = kts.bdry)-1') +
25    geom_line(data=data_fit, color="red", alpha=.10, aes(group=simulation)) +
26    geom_point()
```
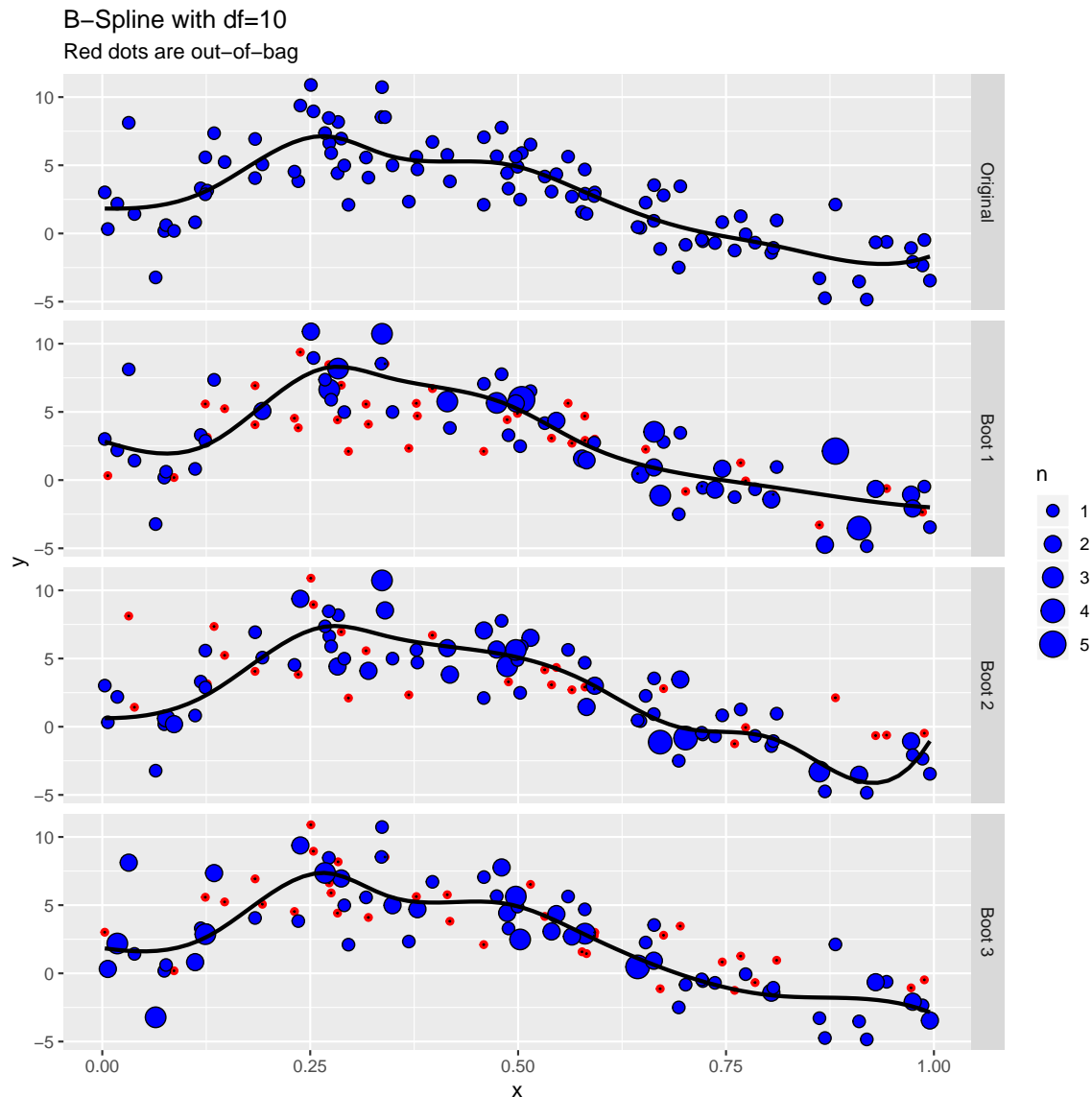


# 4 More Bagging

## 4.1 Out-of-Bag Samples

**Your Turn #1 : Observations not in bootstrap sample**

What is the expected number of observations that will *not* be in a bootstrap sample? Suppose $n$

observations.

Let's look at a few bootstrap fits:

B–Spline with df=10
Red dots are out–of–bag

- Notice that each bootstrap sample does not include about 37% of the original observations.

- These are called *out-of-bag* samples and can be used to assess model fit

  - The out-of-bag observations were not used to estimate the model parameters, so will be sensitive to over/under fitting

- Below, we evaluate the oob error over the spline complexity (df = number of estimated coefficients)
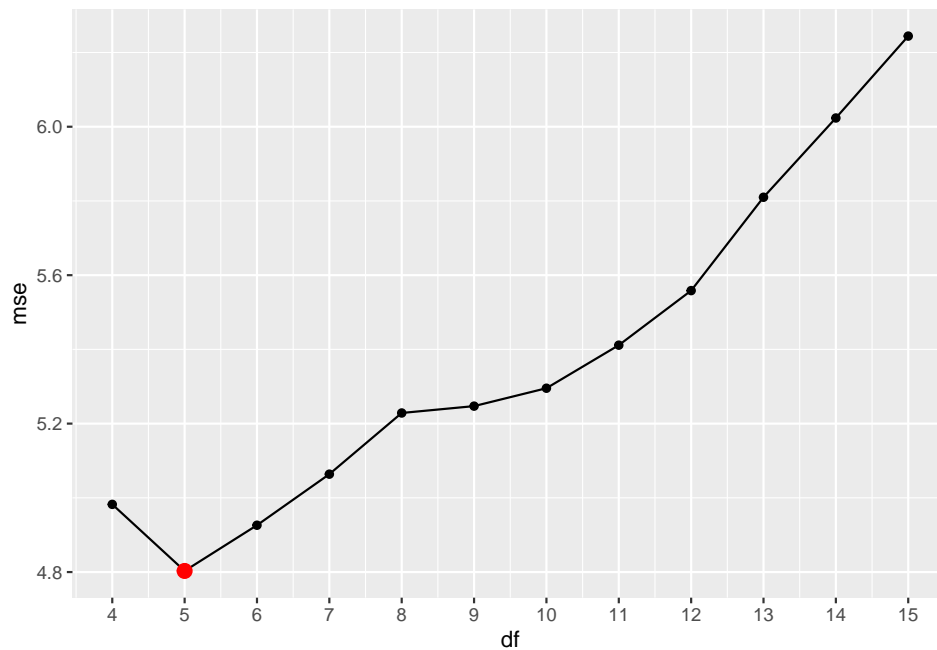
```r
M = 500                          # number of bootstrap samples
DF = seq(4, 15, by=1)            # edfs for spline
results = tibble()
set.seed(2019)

#-- Spline Settings
for(m in 1:M){
  #- sample from empirical distribution
  ind = sample(n, replace=TRUE)        # sample indices with replacement
  oob.ind = which(!(1:n %in% ind))     # out-of-bag samples
  #- fit bspline models
```

```
12    for(df in DF){
13      if(length(oob.ind) < 1) next
14      #- fit with bootstrap data
15      m_boot = lm(y~bs(x, df=df, Boundary.knots=kts.bdry)-1,
16                     data=data_train[ind,])
17      #- predict on oob data
18      yhat.oob = predict(m_boot, newdata=data_train[oob.ind, ])
19      #- get errors
20      sse = sum( (data_train$y[oob.ind] - yhat.oob)^2 )
21      n.oob = length(oob.ind)
22      #- save results
23      results = bind_rows(results,
24                          tibble(m, df, sse, n.oob))
25    }
26  }
27
28  results %>% group_by(df) %>% summarize(mse = sum(sse)/sum(n.oob)) %>%
29    ggplot(aes(df, mse)) + geom_point() + geom_line() +
30    geom_point(data=. %>% filter(mse==min(mse)), color="red", size=3) +
31    scale_x_continuous(breaks=1:20)
```
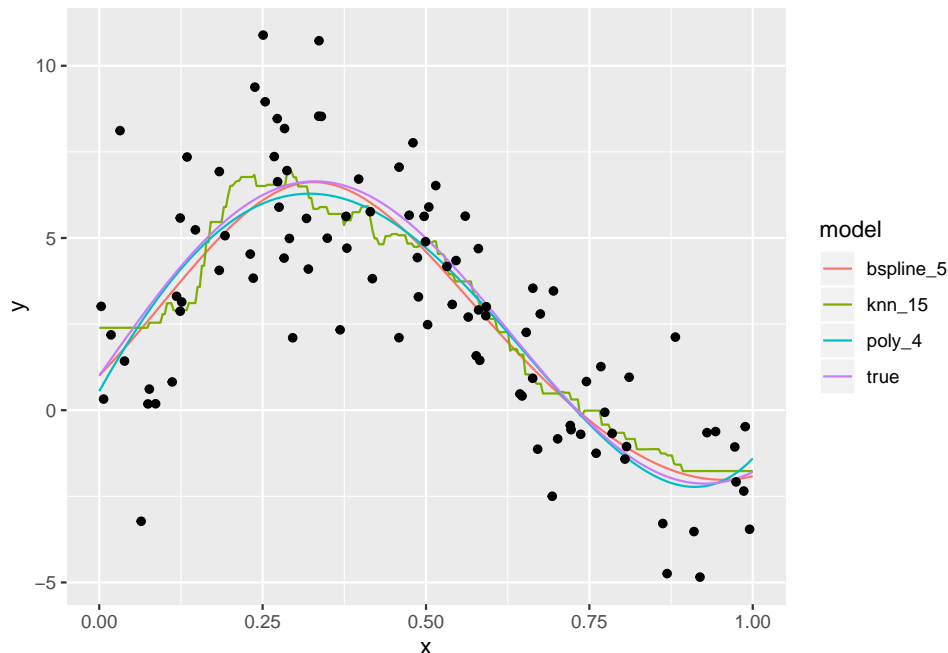


- The minimum out-of-bag error occurs at df=5. This matches the complexity in a polynomial fit from the previous lecture notes.

```
#> Warning: attributes are not identical across measure variables;
#> they will be dropped
```

## 4.2 Number of Bootstrap Simulations

Hesterberg recommends using $M \geq 15{,}000$ for real applications to remove most of the Monte Carlo variability.

- For the examples in class I used much less to demonstrate the principles.

# 5 More Resources

- ESL 5; 8.1-8.4

- ISL 5.2; 7

- The `boot` package and `boot()` function provides some more advanced options for bootstrapping

- What Teachers Should Know About the Bootstrap: Resampling in the Undergraduate Statistics Curriculum, by Tim C. Hesterberg

- R's `tidymodels` package

    - `rsample` for resampling
    - `yardstick` for evaluation metrics
    - `broom` for extracting properties (e.g., estimated parameters) of fitted models in a tidy form

## 5.1 Variations of the Bootstrap

- We have discussed only one type of bootstrap, *nonparametric/empirical/ordinary* where the observations are resampled

- Another option is to simulate from the *fitted model*. This is called the *parametric* bootstrap.

    - For example, in the regression setting, estimate $\hat{\theta}$ and $\hat{\sigma}$
    - Then given the original $X$'s simulate new $y_i^* \mid x_i \sim f(x_i; \hat{\theta}) + \epsilon(\hat{\sigma})$