

Rcpp Portfolio

Rachel Wood

2023-03-27

For this portfolio, we use Rcpp to fit an adaptive kernel smoothing regression model.

We first generate data according to the model

$$y_i = \sin(\alpha \pi x^3) + z_i \quad \text{with} \quad z_i \sim \mathcal{N}(0, \sigma^2)$$

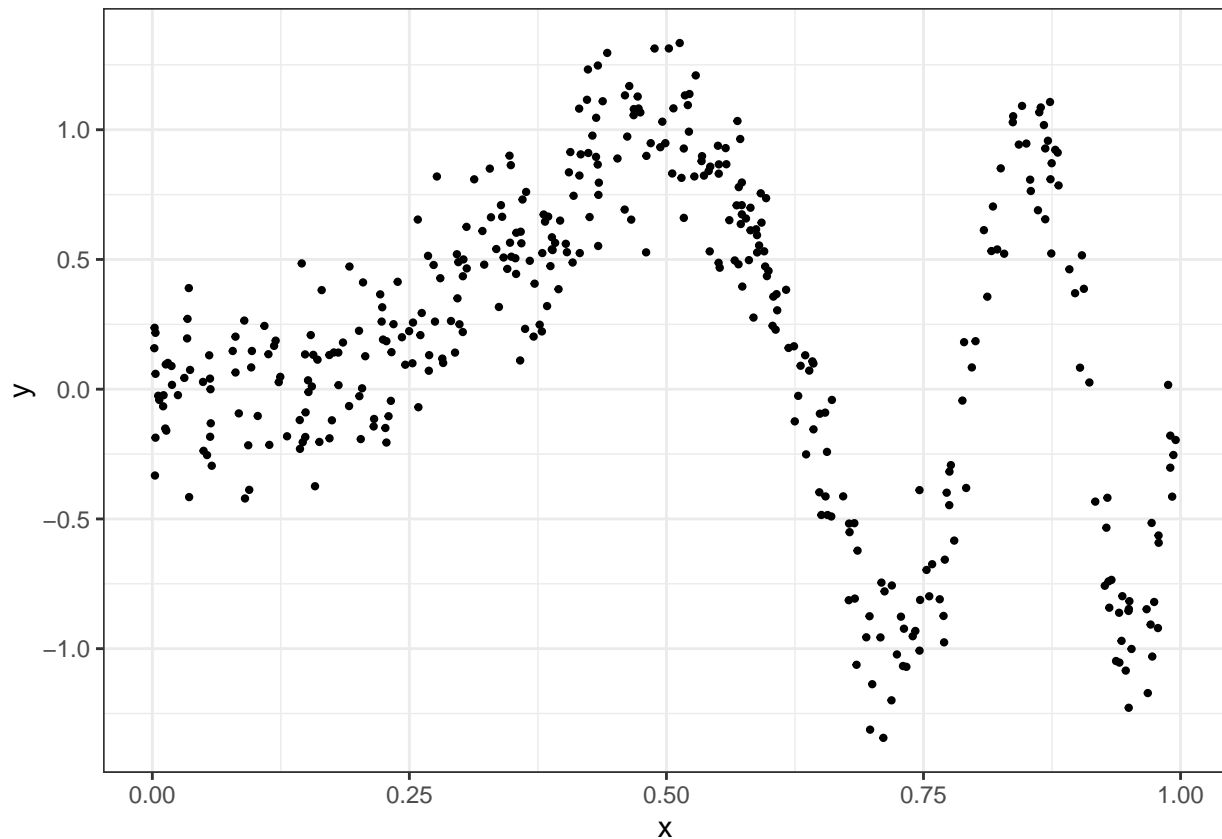
In this case we take $\alpha = 4$ and $\sigma = 0.2$.

```
library(dplyr)
library(ggplot2)
n <- 400
alpha <- 4
sigma <- 0.2

x <- runif(n)
y <- sin(alpha * pi * x^3) + rnorm(n, sd = sigma)

data <- tibble(x = x, y = y)

ggplot(data = data, aes(x, y)) +
  geom_point(size = 0.8)
```



The Kernel Smoother

We model $\mu(x) = \mathbb{E}(y|x)$ by

$$\hat{\mu}(x) = \frac{\sum_{i=1}^n \kappa_{\lambda}(x, x_i) y_i}{\sum_{i=1}^n \kappa_{\lambda}(x, x_i)}$$

where we take κ_{λ} to be a Gaussian kernel with variance λ^2 .

We implement this with the following function:

```
meanKRS <- function(x, y, xnew, lambda){
  n <- length(x)
  nnew <- length(xnew)

  mu <- numeric(nnew)

  for (i in 1:nnew){
    mu[i] <- sum(dnorm(x, xnew[i], lambda)*y) / sum(dnorm(x, xnew[i], lambda))
  }

  return(mu)
}
```

We can now compare the fits for different values of λ :

```
library(tidyr)
xnew <- seq(0, 1, length.out = 1000)
```

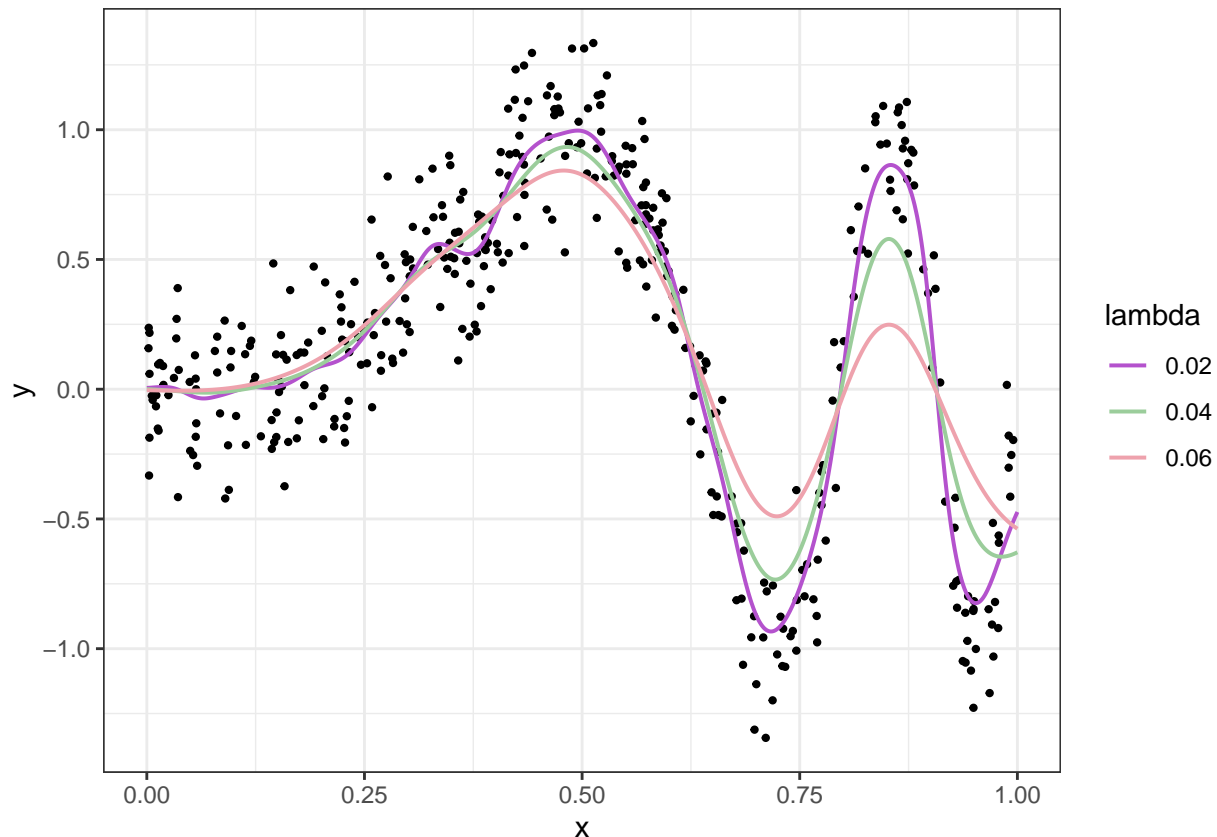
```

smooth_large <- meanKRS(x, y, xnew, lambda = 0.06)
smooth_medium <- meanKRS(x, y, xnew, lambda = 0.04)
smooth_small <- meanKRS(x, y, xnew, lambda = 0.02)

plot_data <- tibble(x = xnew) %>%
  mutate("0.06" = smooth_large,
         "0.04" = smooth_medium,
         "0.02" = smooth_small) %>%
  pivot_longer(cols = c("0.06", "0.04", "0.02"),
               names_to = "lambda",
               values_to = "fitted") %>%
  mutate(lambda = as.factor(lambda))

ggplot() +
  geom_point(data = data,
            aes(x, y), size = 0.8) +
  geom_line(data = plot_data,
           aes(x, fitted, color = lambda), linewidth = 0.7)

```



We

now use Rcpp to write a C++ version of meanKRS():

```

library(Rcpp)

#include <Rcpp.h>
#include <Rmath.h>
using namespace Rcpp;

```

```
// [[Rcpp::export]]

NumericVector meanKRS_Rcpp(const NumericVector x, const NumericVector y, const
↪ NumericVector xnew, const double lambda) {
  int n = x.size();
  int nnew = xnew.size();

  NumericVector mu(nnew);

  for (int i = 0; i < nnew; i++){
    mu[i] = sum(dnorm(x,xnew[i], lambda)*y) / sum(dnorm(x,xnew[i], lambda));
  }

  return mu;
}
```

We check that this function produces the same output as the R version,

```
max(meanKRS(x, y, xnew, lambda = 0.06) - meanKRS_Rcpp(x, y, xnew, lambda = 0.06))
```

```
## [1] 1.221245e-15
```

and compare the performance of the two functions using the `microbenchmark()` function:

```
library(microbenchmark)
microbenchmark("R" = meanKRS(x, y, xnew, lambda = 0.06),
               "Rcpp" = meanKRS_Rcpp(x, y, xnew, lambda = 0.06))
```

```
## Unit: milliseconds
## expr      min       lq      mean   median      uq      max neval
##   R 18.89439 19.16753 19.80777 19.63680 20.02191 22.20248   100
##  Rcpp 12.41657 12.61130 12.75476 12.69977 12.87965 13.56912   100
```

Cross-Validation

We now implement a cross-validation procedure for finding the optimal λ , using the mean squared error of the test set as the metric for determining the fit of our model. We first write the R version of this function:

```
mse_lambda <- function(log_lambda, x, y, x_new, y_new){
  lambda <- exp(log_lambda)

  fitted <- meanKRS(x, y, x_new, lambda)
  return(sum((fitted - y_new)^2))
}

lambda_cv <- function(x, y, groups){
  n <- length(x)

  lambdas <- numeric(nfolds)
  mse <- numeric(nfolds)

  for (i in 1:nfolds){
    x_train <- x[groups != i]
    y_train <- y[groups != i]
```

```

x_test <- x[groups == i]
y_test <- y[groups == i]

solution <- optim(par = 0.02, fn = mse_lambda, x = x_train, y = y_train, x_new =
↪ x_test, y_new = y_test, method = "BFGS")
lambdas[i] <- exp(solution$par)
mse[i] <- solution$value

}

min_ind <- which.min(mse)
return(lambdas[min_ind])
}

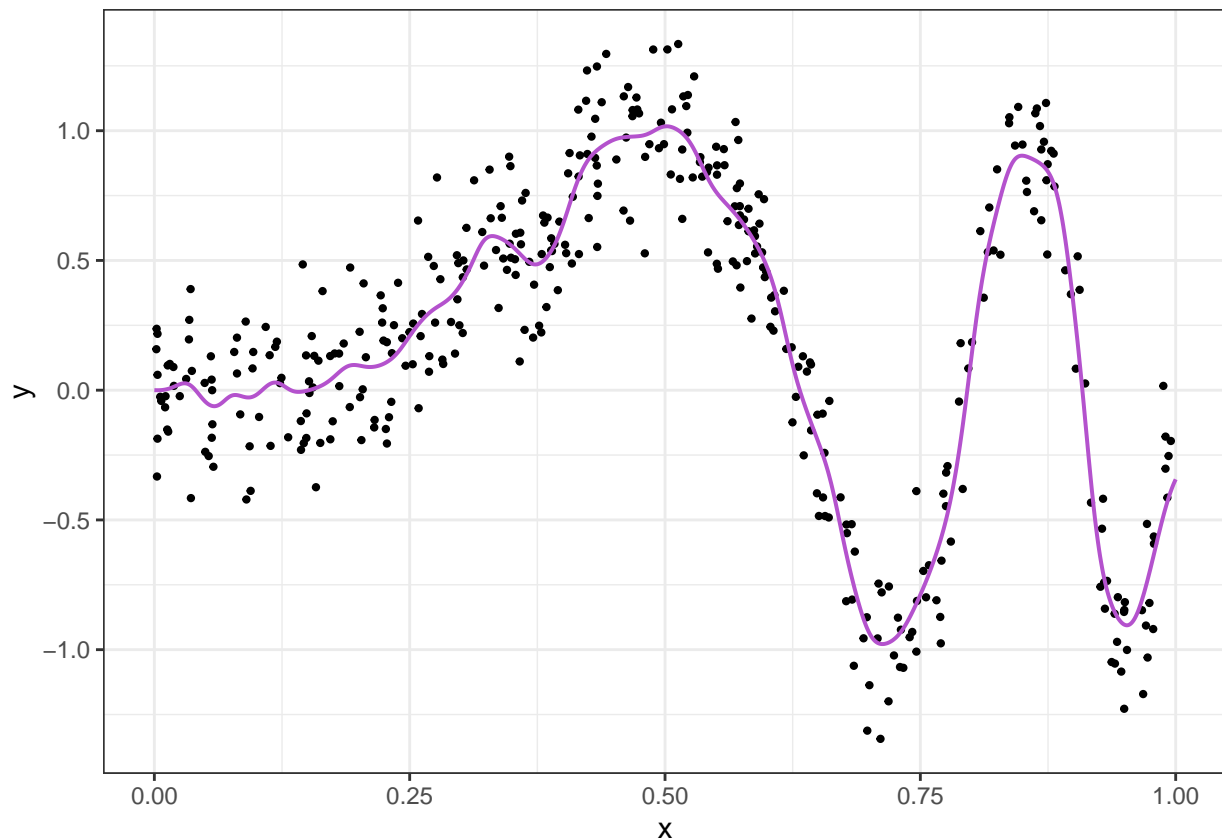
```

We now plot the smooth for the returned value λ to see if this seems reasonable:

```

nfolds <- 5
groups <- sample(rep(1:nfolds, length.out = n), size = n)
hat_lambda <- lambda_cv(x, y, groups)
opt_smooth <- meanKRS(x, y, xnew, hat_lambda)
opt_data <- tibble(xnew = xnew,
  fitted = opt_smooth)
ggplot() +
  geom_point(data = data, aes(x, y), size = 0.8) +
  geom_line(data = opt_data, aes(x = xnew, y = fitted),
    color = "mediumorchid3", linewidth = 0.7)

```



We now write the equivalent function in Rcpp using the `roptim` package:

```
// [[Rcpp::depends(RcppArmadillo)]]
// [[Rcpp::depends(roptim)]]

#include <cmath>
#include <cstdint>

#include <algorithm>

#include <RcppArmadilloExtensions/sample.h>
#include <RcppArmadillo.h>
#include <roptim.h>
#include <functional>
using namespace Rcpp;
using namespace arma;
using namespace roptim;

double mse_lambda(const double lambda, const NumericVector x, const NumericVector y, const
↳ NumericVector xnew, const NumericVector y_new){
    int n = x.size();
    int nnew = xnew.size();

    NumericVector fitted(nnew);

    for (int i = 0; i < nnew; i++){
        fitted[i] = sum(dnorm(x, xnew[i], lambda)*y) / sum(dnorm(x, xnew[i], lambda));
    }
    return sum(pow(fitted - y_new, 2));
}

NumericVector x_train, y_train, y_test, x_test;

// [[Rcpp::export]]

double lambda_cv_Rcpp(const NumericVector x, const NumericVector y, const NumericVector
↳ groups) {

    int n = x.size();

    int nfolds = unique(groups).size();

    NumericVector lambdas(nfolds);
    NumericVector mse(nfolds);

    for (int i = 0; i < nfolds; i++){
        x_train = x[groups != (i+1)];
        y_train = y[groups != (i+1)];

        x_test = x[groups == (i+1)];
```

```

    y_test = y[groups == (i+1)];

class Mse : public Functor {
public:
    double operator()(const arma::vec& log_lambda) override {
        double lambda = exp(log_lambda[0]);
        return mse_lambda(lambda, x_train, y_train, x_test, y_test);
    }
};

Mse fun;
Roptim<Mse> opt("BFGS");

arma::vec initial = {0.02};
opt.minimize(fun, initial);

arma::vec par = opt.par();
lambdas[i] = exp(par[0]);
mse[i] = opt.value();

}

int min_ind = which_min(mse);

return lambdas[min_ind];

}

```

We now verify that this is consistent with our R result

```

Rcpp_hat_lambda <- lambda_cv_Rcpp(x, y, groups)
abs(hat_lambda - Rcpp_hat_lambda)

```

```
## [1] 2.926479e-15
```

We can now compare the computational times:

```

microbenchmark("R" = lambda_cv(x, y, groups),
               "Rcpp" = lambda_cv_Rcpp(x, y, groups))

```

```

## Unit: milliseconds
##  expr      min       lq      mean    median      uq      max neval
##    R 220.6138 225.1243 227.6149 227.1779 229.3568 272.8065   100
##  Rcpp 141.9969 144.6592 145.6935 146.0671 146.8075 150.1780   100

```

The results above show the function has been sped up by a factor of 2.

Lambda as a function of x

We can see merely from looking at the plot, the shape of the function changes at approximately $x = 0.5$, and so these two sections of the function will need different values of λ . We address this by modelling $\lambda = \lambda(x)$. We do this by fitting the model as before for a fixed λ (for this we can use the cross-validated value of λ) and consider the residuals r_1, \dots, r_n .

We can then model these under another KRS with the same λ - producing estimates of the absolute values of the residuals $\hat{v}_1, \dots, \hat{v}_n$. We can then fit

$$\hat{\mu}(x) = \frac{\sum_{i=1}^n \kappa_{\lambda_i}(x, x_i) y_i}{\sum_{i=1}^n \kappa_{\lambda_i}(x, x_i)}$$

with $\lambda_i = \lambda \tilde{w}_i$ where $\tilde{w}_i = \frac{n w_i}{\sum_{i=1}^n w_i}$ for $w_i = \hat{v}_i^{-1}$.

We implement this in R with the `mean_var_KRS()` function:

```
mean_var_KRS <- function(y, x, xnew, lambda){
  n <- length(x)
  nnew <- length(xnew)
  mu <- res <- numeric(n)

  out <- madHat <- numeric(nnew)

  for(ii in 1:n){
    mu[ii] <- sum( dnorm(x, x[ii], lambda) * y ) / sum( dnorm(x, x[ii], lambda) )
  }

  resAbs <- abs(y - mu)
  for(ii in 1:nnew){
    madHat[ii] <- sum( dnorm(x, xnew[ii], lambda) * resAbs ) / sum( dnorm(x, xnew[ii],
↵ lambda) )
  }

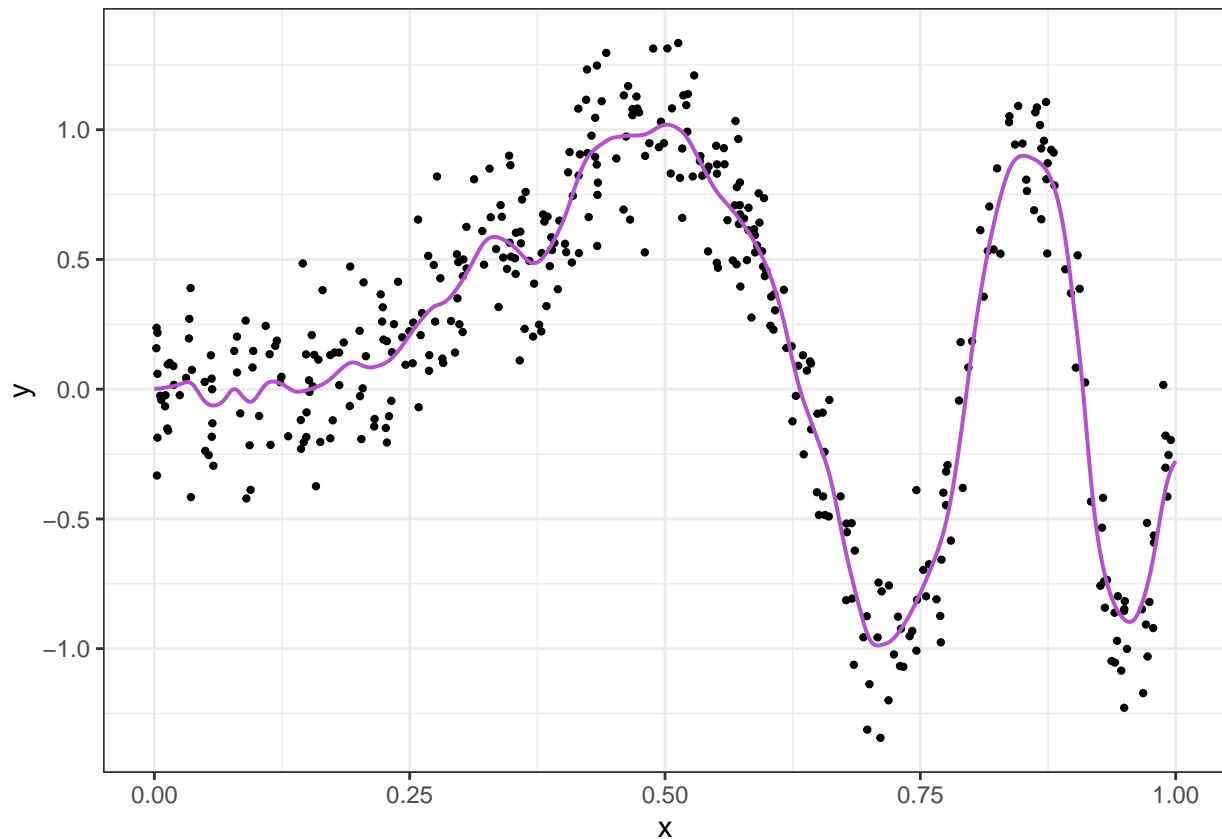
  w <- 1 / madHat
  w <- w / mean(w)

  for(ii in 1:nnew){
    out[ii] <- sum( dnorm(x, xnew[ii], lambda * w[ii]) * y ) /
      sum( dnorm(x, xnew[ii], lambda * w[ii]) )
  }

  return(out)
}
```

and we use our cross-validated value of lambda to get:

```
muSmoothAdapt <- mean_var_KRS( y, x, xnew, hat_lambda)
adapt_data <- tibble(xnew = xnew, fitted = muSmoothAdapt)
ggplot() +
  geom_point(data = data, aes(x, y), size = 0.8) +
  geom_line(data = adapt_data, aes(x = xnew, y = fitted),
    color = "mediumorchid3", linewidth = 0.7)
```

We now write the complementary Rcpp version:

```
#include <Rcpp.h>
#include <Rmath.h>
using namespace Rcpp;

// [[Rcpp::export]]
NumericVector mean_var_KRS_Rcpp(const NumericVector y, const NumericVector x, const
↳ NumericVector xnew, const double lambda){
  int n = x.size();
  int nnew = xnew.size();

  NumericVector mu(n);
  NumericVector res(n);
  NumericVector out(nnew);
  NumericVector madHat(nnew);

  for (int i = 0; i < n; i++){
    mu[i] = sum(dnorm(x,x[i], lambda)*y) / sum(dnorm(x, x[i], lambda));
  }

  NumericVector resAbs = abs(y-mu);

  for (int i = 0; i < nnew; i++){
    madHat[i] = sum(dnorm(x,xnew[i], lambda)*resAbs) / sum(dnorm(x,xnew[i], lambda));
  }
  NumericVector w = 1/madHat;
```

```

w = w/mean(w);

for (int i =0; i <nnew; i++){
  out[i] = sum(dnorm(x, xnew[i], lambda * w[i])*y) / sum(dnorm(x, xnew[i], lambda
↪ *w[i]));
}

return out;
}

```

We now check these two functions produce the same values:

```

muSmoothAdapt_Rcpp <- mean_var_KRS_Rcpp( y, x, xnew, hat_lambda)
max(abs(muSmoothAdapt - muSmoothAdapt_Rcpp))

```

```
## [1] 1.110223e-15
```

Comparing the time gain from this we see the Rcpp version has about half of the computing time:

```

microbenchmark("R" = mean_var_KRS(x, y, xnew, lambda = hat_lambda),
               "Rcpp" = mean_var_KRS_Rcpp(x, y, xnew, lambda = hat_lambda))

```

```
## Unit: milliseconds
##  expr      min       lq      mean   median      uq      max  neval
##    R 43.94724 44.47237 45.61218 45.16792 46.70248 49.86071   100
##  Rcpp 28.33535 28.59356 28.78686 28.75646 28.90781 29.77170   100

```