

# Response to Comments

Thank you for your time and thoughtful comments. I believe that they have helped to greatly improved, not only the manuscript, but the package it describes as well. Please find below a copy of all the comments along with my responses.

## Reviewer 1

### Comment 1.1

**Bug:** A graph is produced as expected if the code below is cut-and-pasted at the command line . However, if the code is put into a file (say `foo.R`) and `source("foo.R")` is executed at the command line, no graph is produced. (But note that if `plot=TRUE` is added to the call to `partial()`, then a plot is produced as expected.)

```
library(randomForest)
data(airquality)
ozone.rf <- randomForest(Ozone ~ ., data=airquality, na.action=na.omit)
library(pdp)
pd <- partial(ozone.rf, pred.var="Temp")
plotPartial(pd)
```

### Response 1.1

This is simply a side effect of using `lattice` to produce the plots. If sourcing a script, then the plot needs to be explicitly printed (e.g., `print(plotPartial(pd))` or `print(partial(fit, pred.var = "x", plot = TRUE))`). This has been clarified in the package documentation and the manuscript (see **Note:** at the bottom of page 3). The fact that this worked when calling `partial` with `plot = TRUE` was a bug (which is now fixed).

### Comment 1.2

I couldn't get `partial` to work with `gbm` (version 2.1.1) models (although I didn't try very hard). Perhaps this is because the `gbm` maintainer changed the `gbm` interface in a non-backcompatible way. The following (where `gbm.mod` is a model built using the `gbm` package)

```
partial(gbm.mod, pred.var="age")
```

gives

```
Error in paste("Using", n.trees, "trees...\n") :
argument "n.trees" is missing, with no default
```

I also tried

```
partial(gbm.mod, pred.var="age", n.trees=100)
```

but that gives

```
Error in .fun(piece, ...) : unused argument (n.trees = 100)
```

### Response 1.2

This has been fixed in `pdp` (version 0.3.0). The user may still have to supply `n.trees` in the call to `partial` (depending on which version of `gbm` they are using). Also, as suggested in another comment, a `/slowtests` directory is underway and will help prevent future issues like this. Thanks for catching this.

### Comment 1.3

Standard partial matching isn't implemented (e.g., `partial(..., type="r")`).

### Response 1.3

Partial matching is now implemented (as of `pdp` version 0.3.0), but only for the `type` argument. Thanks for the suggestion.

### Comment 1.4

The following code incorrectly gives an error message

```
"Error in names(pd_df) <- c(pred.var, "y") :  
  'names' attribute [2] must be the same length as the vector [1]"
```

It seems to work ok if you change the variable names.

```
library(pdp)  
data <- data.frame(V1 = 1:10, V2 = 1:10)  
mod <- lm(V2 ~ V1, data = data)  
partial(mod, pred.var = "V1") # error
```

### Response 1.4

This has been fixed in `pdp` (version 0.3.0). This was simply a side effect of `partial`'s internal use of `plyr::adply` which uses default variable names (e.g., `V1`, `V2`, etc.). This caused the column `V1` in the output data frame to be overwritten, producing the error.

```
library(pdp)  
data <- data.frame(V1 = 1:10, V2 = 1:10)  
mod <- lm(V2 ~ V1, data = data)  
partial(mod, pred.var = "V1") # no more error
```

```
##      V1 yhat  
## 1     1     1  
## 2     2     2  
## 3     3     3  
## 4     4     4  
## 5     5     5  
## 6     6     6  
## 7     7     7  
## 8     8     8  
## 9     9     9  
## 10    10    10
```

### Comment 1.5

I think you should mention in the documentation something about the assumptions you make about accessing the model data. The following (somewhat contrived) example fails with "Error in train[[x]] : subscript out of bounds":

```
data(trees)
foo <- function(data) {
  lm(Volume ~ ., data = data)
}
mod <- foo(trees)
library(pdp)
partial(mod, pred.var = "Girth", plot = TRUE) # fails
```

and the following silently gives misleading results:

```
library(pdp)
data(trees)
mod <- lm(Volume ~ ., data = trees)
partial(mod, pred.var = "Girth", plot = TRUE) # ok
trees <- trees[1:2, ]
partial(mod, pred.var = "Girth", plot = TRUE) # wrong xlim
```

I'm not necessarily saying that you need to resolve the above two issues, but at least mention in the documentation that this kind of thing can happen, and how to avoid it (basically call partial in the same environment used to build the model, and don't change any of the data used to build the model).

### Response 1.5

Agreed! Both the manuscript (see page 3) and the help page for `partial` have been expanded to discuss this. Unfortunately, some functions don't store a copy of the training data. However, I have expanded the code to work harder—especially for S4 methods—to inform the user to supply the training data when necessary. Re-running your first example should give an informative error message:

```
data(trees)
foo <- function(data)
{
  lm(Volume ~ ., data = data)
}
mod <- foo(trees)
library(pdp)
partial(mod, pred.var = "Girth", plot = TRUE) # fails with better error mssg
```

```
## Error in getTrainingData.default(object): The training data could not be extracted from object. Please
```

### Comment 1.6

The following code gives an obscure error message:

```
"Error in seq.default(from = min(train[[x]], na.rm = TRUE),
  to = max(train[[x]], :
  'from' cannot be NA, NaN or infinite"
```

Some better handholding for the user may be helpful here.

```
library(pdp)
data(trees)
mod <- lm(Volume~., data=trees)
partial(mod, pred.var="nosuchvariable", plot=TRUE) # obscure err msg
```

### Response 1.6

This has been fixed in `pdp` (version 0.3.0) by adding an informative error message listing which variables are the cause of the issue.

```
library(pdp)
data(trees)
mod <- lm(Volume~., data=trees)
partial(mod, pred.var=c("Girth", "foo", "bar"), plot=TRUE)
```

```
## Error in partial.default(mod, pred.var = c("Girth", "foo", "bar"), plot = TRUE): foo, bar not found
```

### Comment 1.7

On page 2 of the paper it says “The columns of the data frame are labeled in the same order as the features supplied to `pred.var`, and the last column is always labeled `y`”. What happens if `y` is one of the predictors?

### Response 1.7

The last column is now labeled `yhat`, which is less likely to be a predictor name. When one of the predictors is labeled “`y`” `randomForest::partialPlot` fails while `gbm::plot.gbm` does not. I think regardless of what the default column name is, there is always the possibility of this problem. This has been filed as an issue on the `pdp` GitHub page: <https://github.com/bgreenwell/pdp/issues/22>.

### Comment 1.8

Hastie et al. in the book cited in the paper say “Although such a collection [of partial dependence plots] can seldom provide a comprehensive depiction of the approximation, it can often produce helpful clues”. Therefore in the abstract for the paper, to avoid overselling consider changing “relationship between the outcome and predictors of interest can be easily understood” to “relationship between the outcome and predictors of interest can be more easily understood”.

### Response 1.8

Fixed. Thanks for the suggestion.

### Comment 1.9

On page 3 in the PDF of the paper, add a newline to the error message in the paper:

```
Error: The training data could not be extracted from object. Please
supply the raw training data using the ...
```

### Response 1.9

Fixed. Thanks for pointing that out.

### Comment 1.10

Is there any reason the package is called `pdp` but the main function is called `partial`? Consider giving the function the same name as the package, so e.g., the following works (to help the user get started)

```
library(pdp)
?pdp # fails
```

### Response 1.10

The package was originally called `partial`, but at the time of its original submission this conflicted with a CRAN package of the same name (but with mixed lowercase/uppercase letters). So, I decided to change the name of the package to `pdp`; besides, I think `partial` and `plotPartial` work well together as function names. However, I did add a help page describing the package, main functions, etc. So, `?pdp` no longer fails. However, I am not opposed to creating an alias to `partial` called `pdp` in the future.

### Comment 1.11

`partial` should return the same value regardless of whether or not the `plot` argument is used?

### Response 1.11

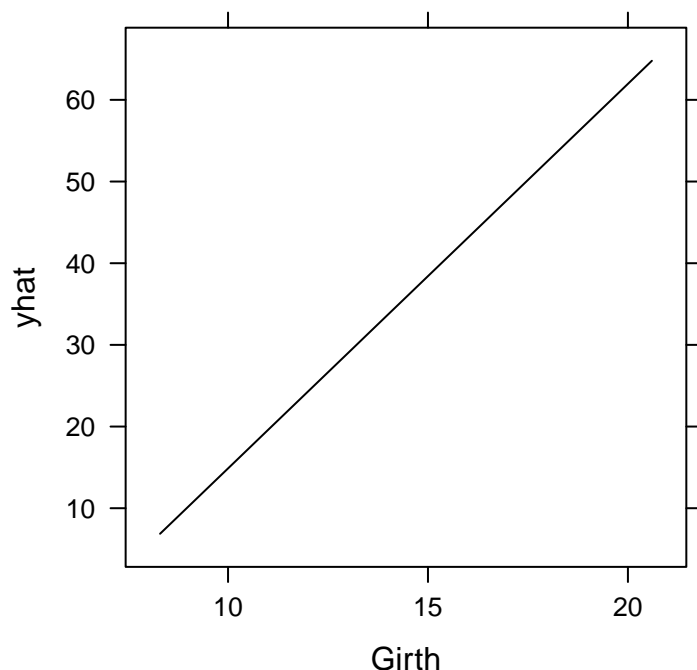
I don't think so, because I want `partial` to behave like a `lattice` function when using `plot = TRUE`; that is, I want it to return a "trellis" object. However, I did add the partial dependence data as an attribute to the returned "trellis" object so that it can easily be extracted if necessary (e.g., if they do not like the plot and don't want to re-run `partial`). For example,

```
library(pdp)
data(trees)
mod <- lm(Volume ~ ., data = trees)
pd.Girth <- partial(mod, pred.var = "Girth") # data only
pdp.Girth <- partial(mod, pred.var = "Girth", plot = TRUE) # plot
identical(pd.Girth, attr(pdp.Girth, "partial.data")) # should be identical
```

```
## [1] TRUE
```

Also, the `lattice` function `trellis.last.object` is now exported when `pdp` is loaded. This would be useful if the user did not store the PDP but needs to retrieve the output data to avoid re-running `partial` (e.g., when computation time is a concern). For example,

```
library(pdp)
data(trees)
mod <- lm(Volume ~ ., data = trees)
partial(mod, pred.var = "Girth", plot = TRUE)
```



```
pdp.Girth <- lattice::trellis.last.object()
pd.Girth <- attr(pdp.Girth, "partial.data")
head(pd.Girth)
```

```
##      Girth      yhat
## 1  8.300000  6.873167
## 2  8.773077  9.100489
## 3  9.246154 11.327811
## 4  9.719231 13.555133
## 5 10.192308 15.782455
## 6 10.665385 18.009777
```

This has been noted in the package documentation (see `?partial`). Thanks for the idea!

### Comment 1.12

A call to say `par(mfrow = c(2, 2))` before calling `plotPartial` is ignored by `plotPartial`. This is a pity because it would allow multiple plots to be put on one page like this:

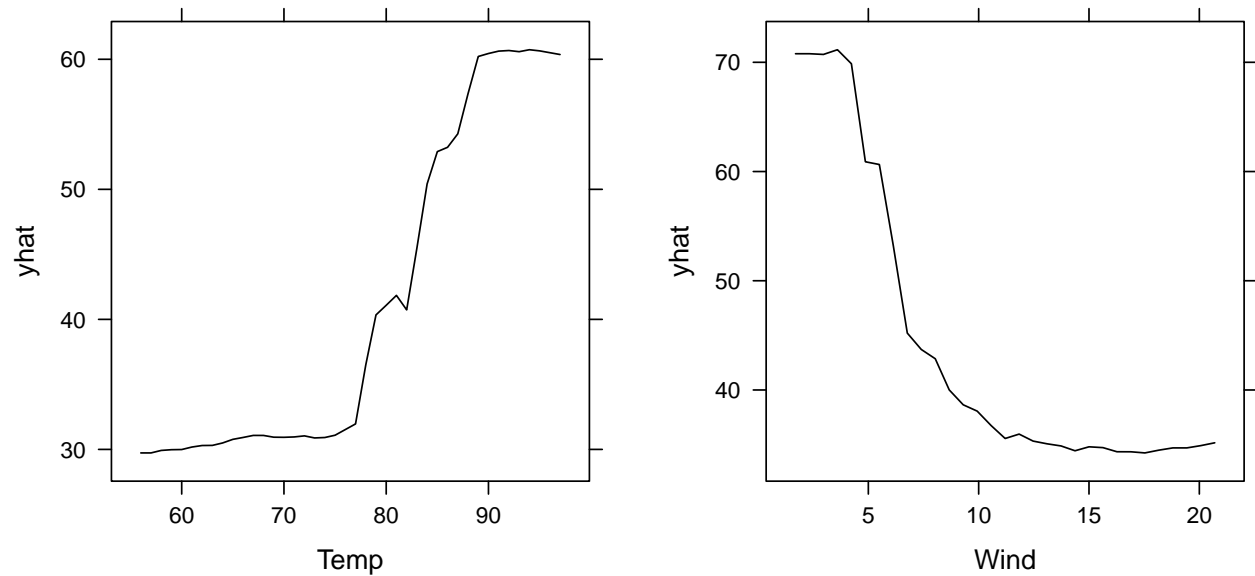
```
par(mfrow = c(1, 2))
partial(ozone.rf, pred.var = "Temp", plot = TRUE)
partial(ozone.rf, pred.var = "Wind", plot = TRUE)
```

### Response 1.12

The command `par(mfrow = c(2, 2))` does not work with "trellis" objects, such as those produced by `lattice`—which `partial` relies on for producing its plots. Other methods are available. For convenience, `pdp` now imports the `grid.arrange` function from the `gridextra` package (this has been discussed in the manuscript and incorporated into the examples where appropriate).

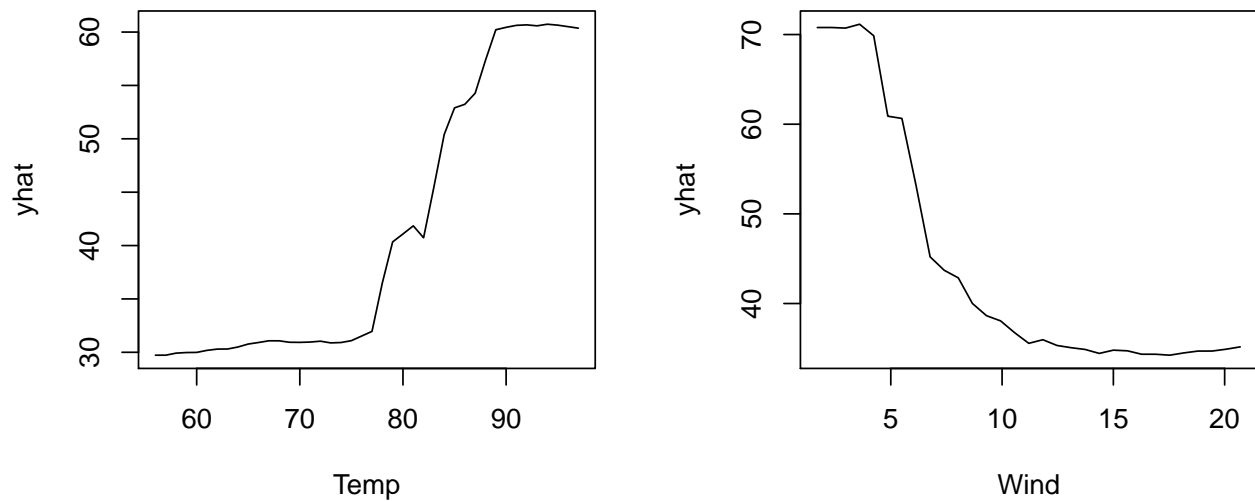
```
grid.arrange(
  partial(ozone.rf, pred.var = "Temp", plot = TRUE),
```

```
partial(ozone.rf, pred.var = "Wind", plot = TRUE),
ncol = 2
)
```



However, since `partial` returns a data frame by default, using `par` would work in the following case

```
par(mfrow = c(1, 2))
plot(partial(ozone.rf, pred.var = "Temp"), type = "l")
plot(partial(ozone.rf, pred.var = "Wind"), type = "l")
```



### Comment 1.13

It would be nice if the functions automatically plotted all the (important) variables in a grid of plots, without forcing the user to explicitly specify them.

### Response 1.13

I agree, this would be a nice feature, but the problems I see are (1) not all models naturally emit a “variable importance” score; the `caret` package has come a long way in this regard, and (2) too many plots will look cluttered. At this point, I’d rather encourage the user to select which variables they would like to see plotted. I am not opposed to including this feature in the future. I added this as a feature request: <https://github.com/bgreenwell/pdp/issues/19>. Thanks for the suggestion!

### Comment 1.14

It would be nice if the functions automatically got the data from the model used to build the model in more cases without forcing the user to explicitly specify the data. Since XGBoost seems now quite popular, perhaps the authors of that package could be persuaded to make the original model data available for functions like `partial` (currently XGBoost uses a custom object `xgb.DMatrix` which is currently write-only—one can’t access the data stored in `xgb.DMatrix`).

### Response 1.14

`kernlab::ksvm` does not store a copy of the training data either, as far as I can tell. As time goes on, the coverage will certainly improve, which will make it easier on the user. For example, `partial` now works with `caret` which, by default, stores the training data in a component called `trainingData`. Also, I suspect that the package maintainers for `xgboost` (and similar packages designed for working with “large” data sets) would be hesitant in storing a copy of the training data with the fitted model for memory reasons.

### Comment 1.15

Consider adding a `slowtests` directory (below `pdp/inst/` or elsewhere) that has tests for all the models supported by the package (Table 1 in the paper) to easily check back compatibility when packages change. (For justification for such tests see for example my comments on `gbm` back compatibility above.) These tests will probably be quite slow, but that won’t be an issue because they won’t be invoked by CRAN check; the slow tests would be just for the maintainer of the `pdp` package to run manually.

### Response 1.15

This is a good suggestion and is currently under development: <https://github.com/bgreenwell/pdp/blob/master/slowtests/slowtests.R>. In fact, the `gbm` maintainers have already contacted me about a breaking change caused by the next planned release of their package.

## Reviewer 2

### Comment 2.1

The Introduction immediately introduces the PDP. You must first motivate why it is important. In a linear model, they are not important. In a linear model with simple interactions, likewise. They first become important if you have non-parametric machine-learning techniques. Further, they are not the only idea extant when trying to explain a predictor’s impact on the response. This paper would benefit if you can give a bit more of a lit review here.



## Response 2.1

I expanded the introduction a bit to try and motivate why PDPs are important (the new material links directly to the Boston housing example used throughout the manuscript). I also included more background on the `plotmo` and `ICEbox` packages (and their corresponding methodologies). `plotmo` very recently introduced a `partdep` option for displaying PDPs. Because of this, more detail was added to differentiate `pdp` and motivate the use of `partial`.

I agree that PDPs are not as useful in “simple” linear models (e.g., few variables, only two-way interactions, etc.), but certainly there will be cases where PDPs will be extremely useful (e.g., many terms, complex interactions, etc.). In fact, linear models are still common in machine learning/computer science/data science (they just typically contain many terms). I also agree that PDPs are not the only, or sometimes best, approach to understanding the impact of a subset of the predictor’s on the response. Brief descriptions and references are given for many alternatives: marginal model plots (`plotmo`), individual conditional expectation plots (`ICEbox`)—which are described as an alternative to PDPs when strong interactions are suspected, effect plots (`effects`), and a few other techniques available in the `car` package.

## Important Clarifications Required:

### Comment 2.2

Top of p2. The `ICEbox` package allows for general construction and visualization of PDP’s. However, it is limited as you only get one dimensional representations. The selling point of your package should focus on the flexibility of graphics and the illustration of non-extrapolation zones.

## Response 2.2

`ICEbox` does not actually construct PDPs, rather individual curves for each observation (which can be difficult to interpret when you have millions of observations). I think flexible graphics and non-extrapolation zones are only half of the selling point here. I think that being able to run `partial` in parallel is huge, and also that this package works with many types of machine learning models—not just tree-based models. For example, as described in the manuscript, `pdp` can produce single or multi-predictor PDPs for conditional random forests and XGBoost models (among many others) in parallel. I expanded the last example to describe the interface with `caret`, which I think opens the door for wonderful exploratory data analysis pipelines (e.g., choose model → plot importance → plot relationships, etc.). Also, a new section was added to the manuscript: **User-defined prediction functions**. This section describes the new `pred.fun` argument (as suggested by a later comment) that makes `partial` even more flexible.

### Comment 2.3

All `pdp` figures throughout the paper: Can you make explicit in the caption what code produces the left and right subfigures? I would use the `subfigure` package in latex and then put a comment in the R code e.g. `# Fig 3b`.

## Response 2.3

Rather than using subfigures, which I think would look messy with all the code, I made quite an attempt to comment which pieces of code produced which figure. Also, in response to another reviewer’s comments, the code which produces figures containing multiple plots is better described and demonstrated. See, for example, the **Multi-predictor PDPs** section.

**Edits:****Comment 2.4**

p1 abstract: strict assumptions imposed by traditional statistical models

**Response 2.4**

Added some clarification here; in particular, normality, constant variance, and linearity.

**Comment 2.5**

p1 pairs should be using curly braces of angle braces, not square braces as this is easily confused with closed intervals.

**Response 2.5**

Fixed. Switched to curly braces.

**Comment 2.6**

p2 unless the editor disagrees, there is no reason to add `install.packages("pdp")` in the text of the paper as it is assumed a user knows how to install a package.

**Response 2.6**

Fixed. Code chunk was removed.

**Comment 2.7**

p3 “Error: The training...” should be wrapped.

**Response 2.7**

Fixed. The text is now wrapped appropriately. Thanks for catching this!

**Comment 2.8**

Fig 6: the sidebar should be labeled “ozone”. Also, would it be terrible to label all axes even though it is redundant?

## Response 2.8

The sidebar (i.e., colorkey) has been labeled, but for this example only. As sad as it sounds, labeling the colorkey is not as trivial as we'd expect for `lattice` plots. The code for doing so has been incorporated in the example in question (see pages 8-9). Also, the text describing the first code chunk that produces a display with a colorkey now contains a footnote pointing to the example where the code to produce a label can be found. This is also motivation for why one of the examples showed how to use `ggplot2` for plotting the results from `partial`—`ggplot2` makes it a trivial matter to label the colorkey. I don't think it would be terrible to label all the axes. However, in my opinion, labelling all the axes would (1) looks cluttered, and (2) require more unnecessary code that would probably distract the reader. This is the standard axis labelling scheme used in trellis plots from S-plus to R (via `lattice`).

## Comment 2.9

Fig 7: as beautiful as this is, you should label the sidebar as “logit” or “centered logit” instead of “y”.

## Response 2.9

Fixed. Color key labeled “Centered logit”.

## Suggestions for the software

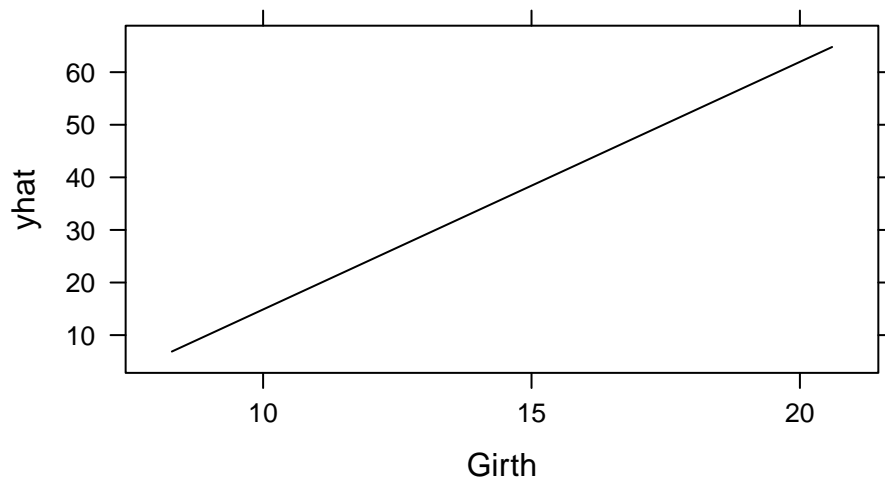
### Comment 2.10

I understand that you like the feature of having a one-shot function partial to both generate and plot the PDP. I realize that the default is `FALSE`. However, I still think it's a disservice to the user as constructing PDP's is computationally intense and if they get the plot wrong or it's not to their liking, they will have to regenerate the whole thing again. So, it is my opinion that you should take it out and force two steps: generation and illustration.

### Response 2.10

While I am not suggesting that the `partial` function needs to mimic the behavior of the PDP functions in `randomForest` and `gbm`, I do find it useful to be able to plot the results without forcing the user to store them first. However, I agree with your concern which is why `partial` has `plot = FALSE` as the default; at the bottom of the second paragraph of section **Constructing PDPs in R** (and in the package documentation), I have added a note emphasizing the use of `plotPartial`. Furthermore, I have made it possible to retrieve the raw PDP data even when `plot = TRUE`. This is because the “`trellis`” object returned by `partial` / `plotPartial` now contains an additional `partial.data` attribute. This is described in the package documentation. For example,

```
library(pdp)
data(trees)
mod <- lm(Volume ~ ., data = trees)
partial(mod, pred.var = "Girth", plot = TRUE)
```



```
pdp.Girth <- lattice::trellis.last.object()
pd.Girth <- attr(pdp.Girth, "partial.data")
head(pd.Girth)
```

```
##      Girth      yhat
## 1  8.300000  6.873167
## 2  8.773077  9.100489
## 3  9.246154 11.327811
## 4  9.719231 13.555133
## 5 10.192308 15.782455
## 6 10.665385 18.009777
```

#### Comment 2.11

The package could be more powerful if you allowed a predict function as a parameter. For instance, you cannot currently use random forests to illustrate classification probabilities.

#### Response 2.11

This is an excellent suggestion. First, I have changed the behavior of `...` to refer to additional arguments in `predict` (e.g., the `n.trees` argument in `predict.gbm`). Second, `partial` gained a new argument, `pred.fun`, where the user can provide an optional prediction function. This makes `partial` far more flexible. For example, you can use this argument to construct PDPs on the probability scale, or to construct ICE curves. This is illustrated in the new section labeled **User-defined prediction functions**.

#### Comment 2.12

p7 I believe there are packages allowing for parallelization in Windows. Can you switch the example so this works on all platforms?

#### Response 2.12

Good suggestion. I tweaked the example to run on both Windows and Unix-like systems. This is also reflected in the manuscript text.