# A Computational Analysis of the Dynamics of R Style
# Based on 94 Million Lines of Code from All CRAN Packages in the Past 20 Years

**Chia-Yi Yen**
R–Ladies Taipei / University of Mannheim
yen.chiayi@gmail.com

**Mia Huai-Wen Chang**
R–Ladies Taipei / Akelius Residential Property AB
mia5419@gmail.com

**Chung–hong Chan**
Hong Kong R User Group / University of Mannheim
chung-hong.chan@mzes.uni-mannheim.de

## ⚠ THE PROBLEM: There are so many programming style variations (PSV). Which one do you use?

```
myFunction<- function(x,y,z = TRUE)
{
  if(z) {
    x + y
}}
```

```
my_function=function(x, y, z = T){;
if(z) { x+y };
};
```
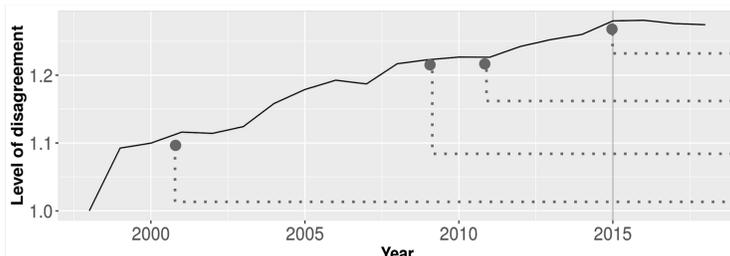
```
Myfunction<-function(x,y,z=T){
    if(z)
x + y }
}
```

```
MYFUNCTION = function(x, y, z=TRUE){
  if(z) {
     x+y
  }
}
```

4 out of at least possible **7168** variations.

❓ Which of the above cannot be correctly evaluated?

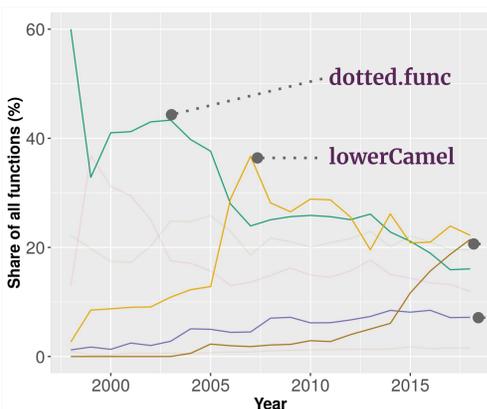## We analysed the distribution of PSV from 1998 to 2018. A consensus is building…



**2015:** *R Packages* the book
**2011:** RStudio / devtools
**2009:** Google Style Guide
**2001:** S4 / = as assignment

# level of disagreement is quantified by normalized Shannon entropy of the PSV distribution per year. Lower value indicates a dominance of certain styles.
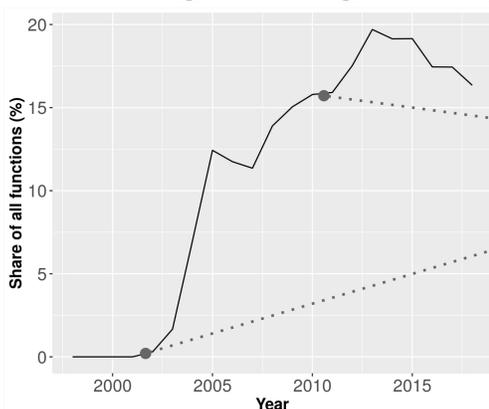
---

# WHAT DRIVES PSV?

### Effect of Style-guides:
naming conventions



- dotted.func
- lowerCamel
- lower__snake (tidyverse)
- UpperCamel (Google)

# de-emphasize 3 naming conventions (ALLUPPER, alllower, other) for simplicity
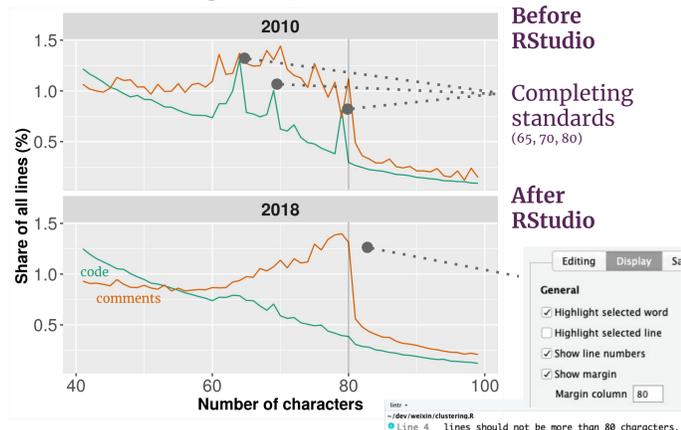
### Effect of Introducing a New Language Feature: using = as assignment operator
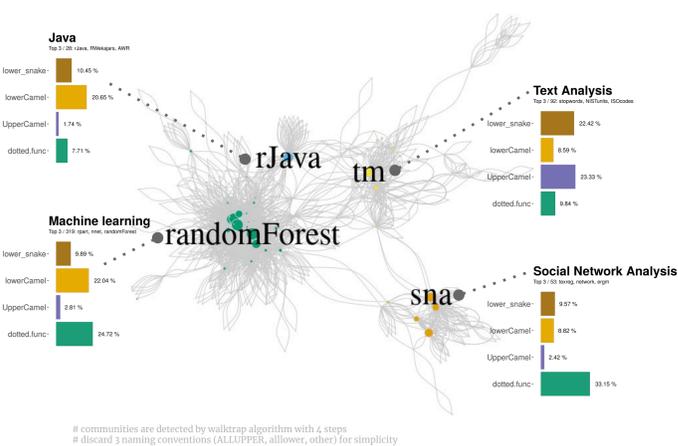


**2010** Ihaka used this style in his paper.
**2001** Introduction of the feature (R 1.4.0)

### Effect of Editors :
line length (options, linters)



**Before RStudio**
Completing standards (65, 70, 80)
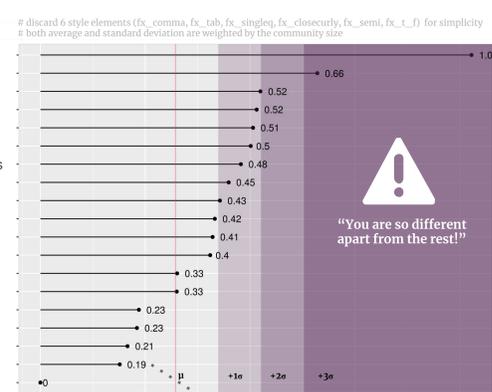
**After RStudio**

---

# COMMUNITY-SPECIFIC VARIATIONS

### CRAN Packages Dependency Network
(subgraph of 4 selected communities and the distribution of naming conventions for each community)



**Java**
Top 3: SB, rJava, RWeekgen, ARR
- lower_snake — 10.45 %
- lowerCamel — 20.65 %
- UpperCamel — 1.74 %
- dotted.func — 7.71 %

**Machine learning**
Top 3: 244, sport, mvt, randomForest
- lower_snake — 9.89 %
- lowerCamel — 22.04 %
- UpperCamel — 2.81 %
- dotted.func — 24.72 %

**Text Analysis**
Top 3: B, stopwords, NLPutils, ISOcodes
- lower_snake — 22.42 %
- lowerCamel — 8.99 %
- UpperCamel — 23.33 %
- dotted.func — 9.94 %

**Social Network Analysis**
Top 3: GGV, txmig, network, expm
- lower_snake — 9.57 %
- lowerCamel — 8.82 %
- UpperCamel — 2.41 %
- dotted.func — 33.15 %

# communities are detected by walktrap algorithm with 4 steps
# discard 3 naming conventions (ALLUPPER, alllower, other) for simplicity

### Divergence in Styles among Communities
(dispute over 10 style-elements among 18 large communities)

| | fx_opencurly | fx_integer | fx_infix | fx_assign | % with this feature |
|---|---|---|---|---|---|
| `func <- function(x)` / `func <- function(x) {` | `x <- 1` / `x <- 1L` | `x <- 3+3` / `x <- 3 + 3` | `x = "Hello World!"` / `x <- "Hello World!"` | | |
| Graphics | 80% | 18% | 18% | 7% | |
| Neuroscience | 55% | 84% | 75% | 9% | |
| Insurance and Actuary | 42% | 76% | 66% | 46% | |
| Java | 26% | 37% | 47% | 17% | |
| Sparse Matrix | 18% | 37% | 26% | 23% | |
| Genetics | 49% | 77% | 71% | 24% | |
| Social Network Analysis | 32% | 74% | | 10% | |
| Finance | 39% | 47% | 39% | 31% | |
| Time, Date, and Money | 24% | 47% | 31% | 2% | |
| Text Analysis | 15% | 39% | 45% | 7% | |
| Numerical Optimization | 36% | 78% | 68% | 12% | |
| RStudio-related | 15% | 46% | 25% | 13% | |
| base | 40% | 74% | 63% | 19% | |
| GPS and Geography | 38% | 70% | 63% | 19% | |
| Machine learning | 40% | 65% | 49% | 27% | |
| RCpp | 18% | | 34% | 15% | |
| Graph data structure | 38% | 55% | 51% | 6% | |
| Image Plotting | 32% | 69% | 56% | 22% | |
| Average | 28% | 61% | 46% | 17% | |

### Who is the "Naughty, Naughty"?
(numbers are the Euclidean distance to average community)



# discard 6 style elements (fx_comma, fx_tab, fx_singleq, fx_closecurly, fx_semi, fx_t_f) for simplicity
# both average and standard deviation are weighted by the community size

- Graphics — 1.0
- Neuroscience — 0.66
- Insurance and Actuary — 0.52
- Java — 0.52
- Sparse Matrix — 0.51
- Genetics — 0.5
- Social Network Analysis — 0.48
- Finance — 0.45
- Time, Date, and Money — 0.43
- Text Analysis — 0.42
- Numerical Optimization — 0.41
- RStudio-related — 0.4
- base — 0.33
- GPS and Geography — 0.33
- Machine learning — 0.23
- RCpp — 0.23
- Graph data structure — 0.21
- Image Plotting — 0.19
- Average — 0

"You are so different apart from the rest!"

$μ = 0.322$, $σ = 0.102$

$d = || \bar{b} - \bar{b} ||$

---

# SUMMARY

### Like Likes Like
(popularity of naming styles among 18 large communities)

>40% using lower__snake



Legend: dotted.func, ALLUPPER, UpperCamel, other, alllower, lowerCamel, lower_snake

0% using lower__snake

### Consensus-based Style
(numbers are % of functions using that style-element)

Use `lowerCamel` (22.2%) or `snake_case` (21.5%)
Use `<-` to assign (83.7%)
Add a space after commas (83.7%)
Don't use `T` / `F` (95.1%)

```
softplusFunc <- function(value, leaky = FALSE) {
    if (leaky) {
        warnings("using leaky RELU!")
        return(ifelse(value > 0, value, value * 0.01))
    }
    return(log(1 + exp(value)))
}
```

Don't use tab to indent (89.2%)
Use `}` on its own line, unless before `else` (89.2%)
Don't explicitly type integers (i.e. `1L`) (60.4%)
Don't terminate lines with `;` (94.6%)
Add spaces around infix operators (55.3%)
Use double quotes for strings (87.6%)
Use same line `{` then a newline (70.0%)