

# Broom: Converting Statistical Models to Tidy Data Frames

David Robinson

6/28/2016

What is tidy data?

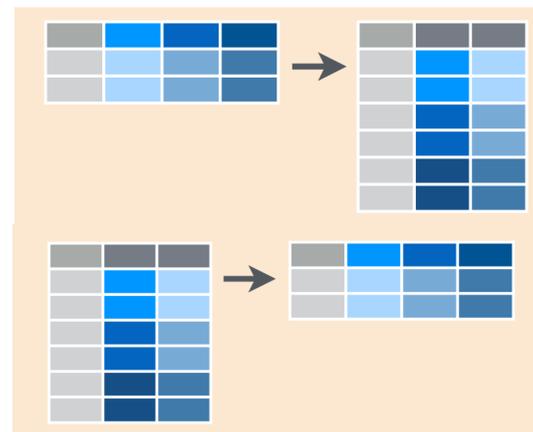
# Data frames arranged as:

- One row for each *observation*
- One column for each *variable*
- One table for each *type of observational unit*

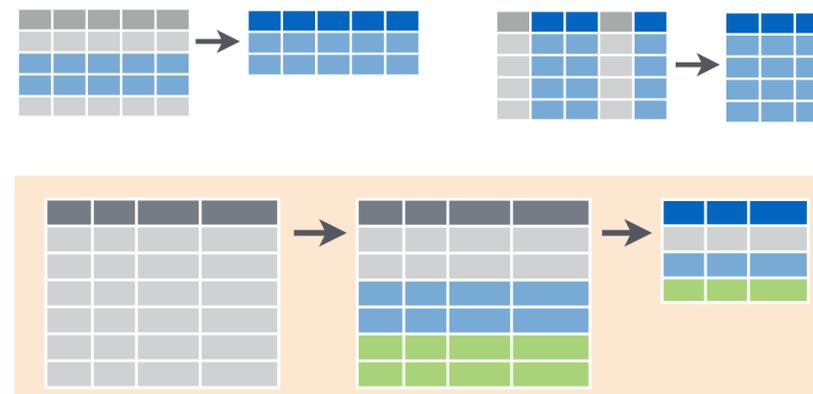
**For details, see [Tidy Data \(Wickham 2014\)](#)**

# “Tidy tools” work with tidy data frames

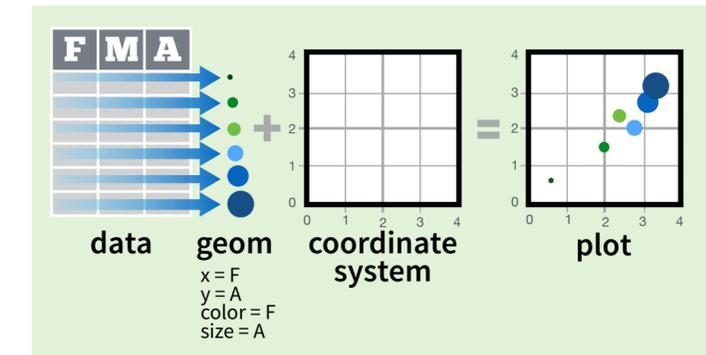
## tidyr



## dplyr



## ggplot2

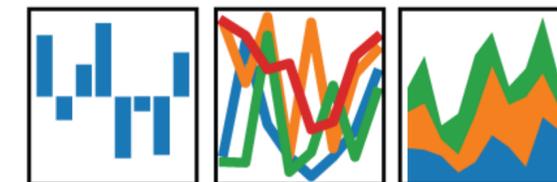


## data.table

| What?  | Example   | Notes   | Output  |
|--|---|---|---|
| Doing <code>sum</code> by group.   | <code>DT[, .(V4.Sum = sum(V4)), by=V1]</code>                     | Calculates the sum of V4, for every group in V1.  | V1 V4.Sum<br>1: 1 36  |
| Doing <code>sum</code> by several groups using <code>.</code> .              | <code>DT[, .(V4.Sum = sum(V4)), by=.(V1, V2)]</code>              | The same as above, but for every group in V1 and V2.                                      | V1 V2 V4.Sum<br>1: 1 A 8<br>2: 2 B 10<br>3: 1 C 12<br>4: 2 A 14<br>5: 1 B 16<br>6: 2 C 18 |
| Call functions in <code>by</code> .  | <code>DT[, .(V4.Sum = sum(V4)), by=sign(V1-1)]</code>             | Calculates the sum of V4, for every group in <code>sign(V1-1)</code> .                    | sign V4.Sum<br>1: 0 36<br>2: 1 42   |
| Assigning new column names in <code>by</code> .                              | <code>DT[, .(V4.Sum = sum(V4)), by=.(V1, O1 = sign(V1-1))]</code> | Same as above, but with a new name for the variable we are grouping by.                   | V1 O1 V4.Sum<br>1: 0 36<br>2: 1 42  |
| Grouping only on a subset by specifying <code>i</code> .                     | <code>DT[1:5, .(V4.Sum = sum(V4)), by=V1]</code>                  | Calculates the sum of V4, for every group in V1, after subsetting on the first five rows. | V1 V4.Sum<br>1: 1 9<br>2: 2 6   |
| Using <code>.N</code> to get the total number of observations of each group. | <code>DT[, .N, by=V1]</code>                                      | Count the number of rows for every group in V1.   | V1 N<br>1: 1 6<br>2: 2 6  |

pandas

$$y_i t = \beta' x_{i t} + \mu_i + \epsilon_{i t}$$



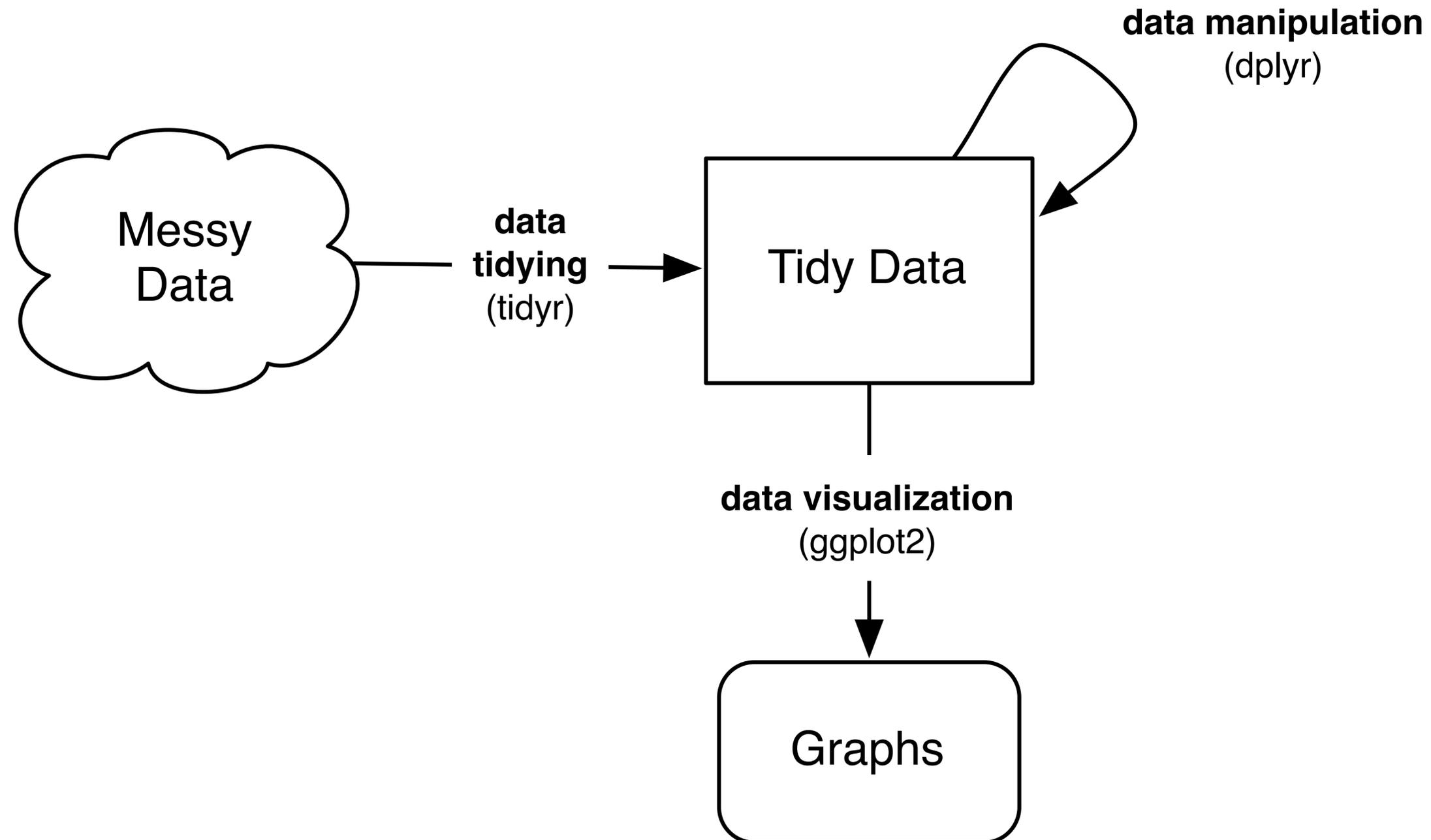
Source: RStudio: Data Wrangling Cheatsheet

RStudio: Data Visualization Cheatsheet

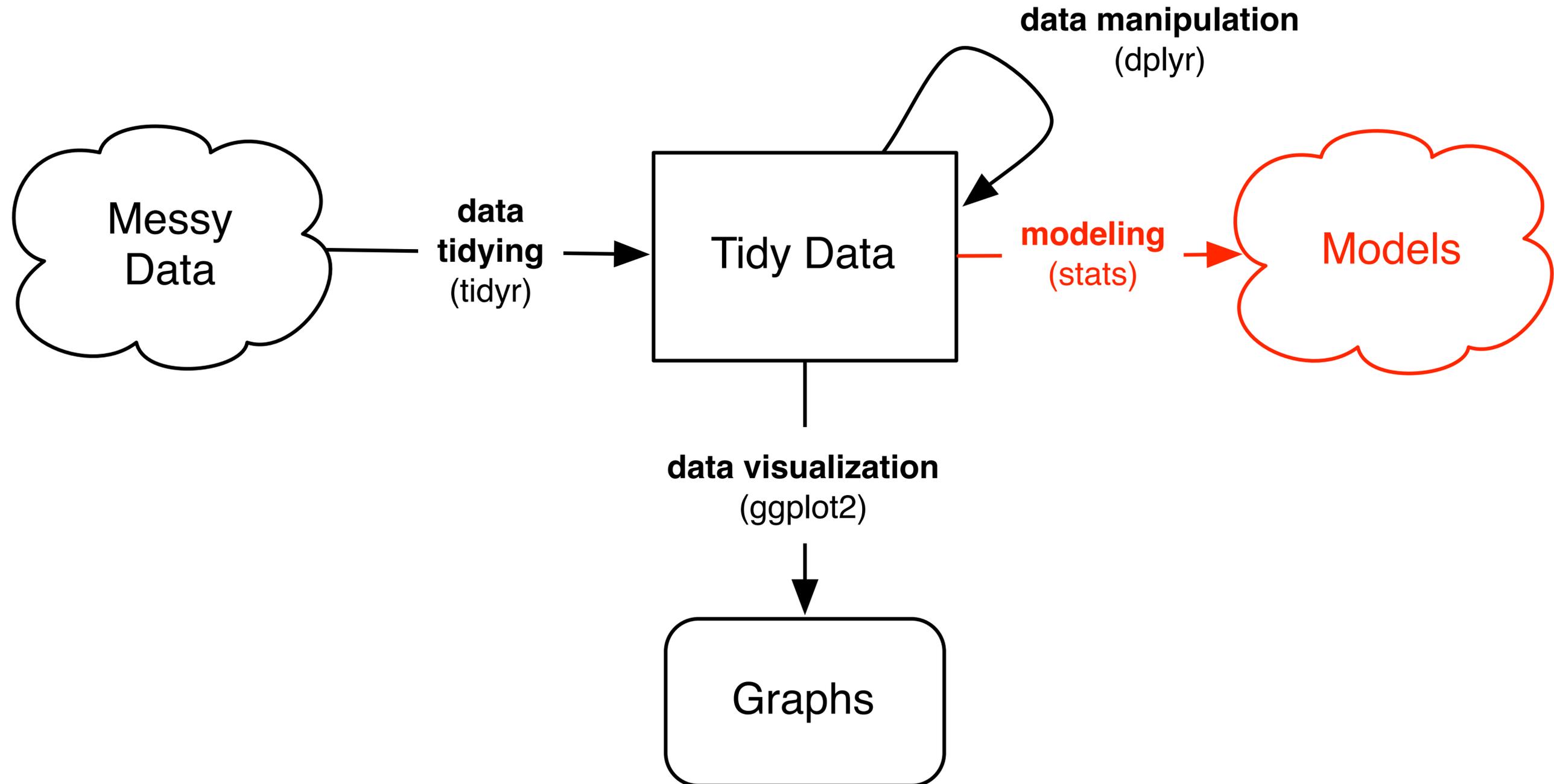
DataCamp: Data Analysis The data.table Way (DataCamp)

<http://pandas.pydata.org/>

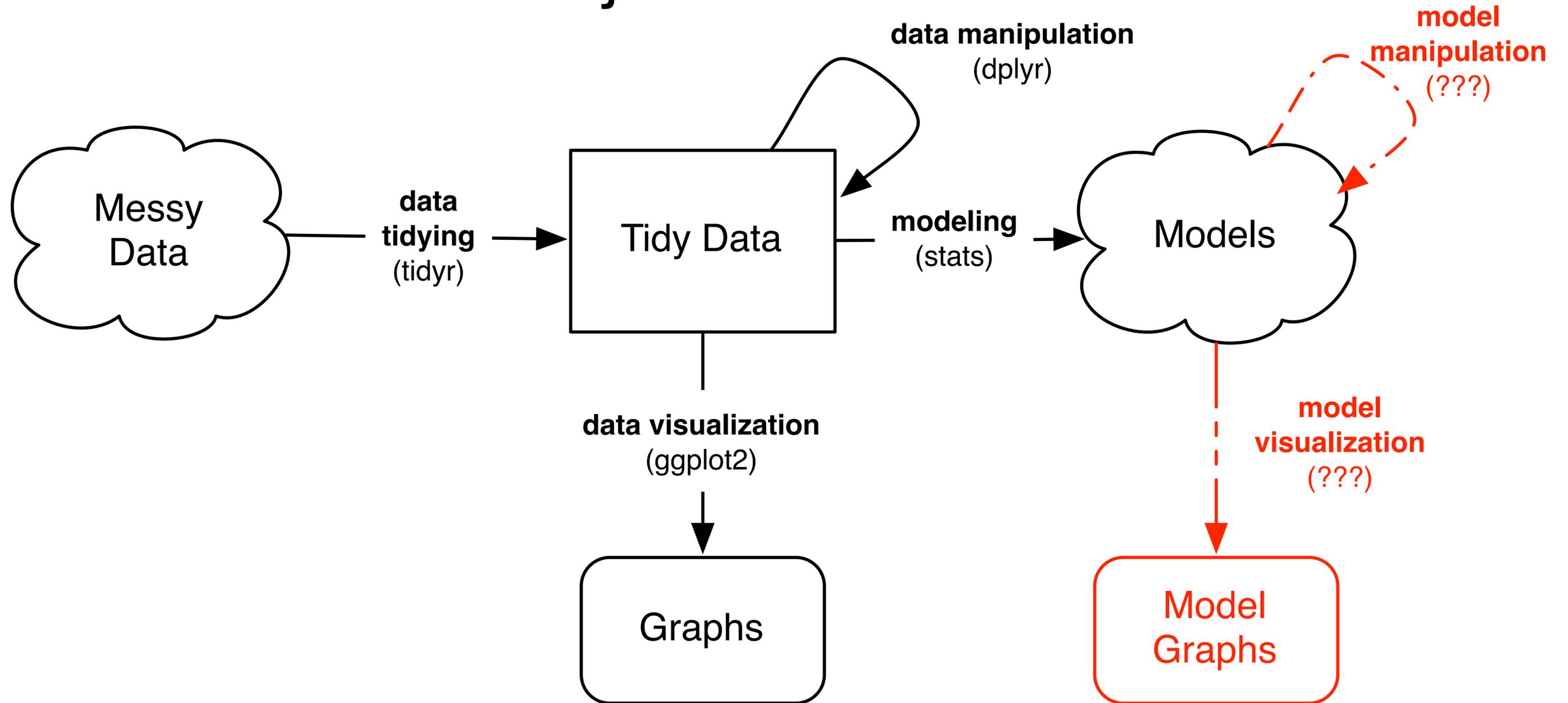
# Tidy tools work together in exploratory data analysis



# Everything works well until...



# Visualizing and manipulating model objects is difficult



**Model objects are  
messy**

Example:  
linear regression

# What's “messy” about a linear regression?

```
> lmfit <- lm(mpg ~ wt + qsec, mtcars)
```

# What's "messy" about a linear regression?

```
> summary(lmfit)
```

```
Call:
```

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

```
Residuals:
```

| Min     | 1Q      | Median  | 3Q     | Max    |
|---------|---------|---------|--------|--------|
| -4.3962 | -2.1431 | -0.2129 | 1.4915 | 5.7486 |

```
Coefficients:
```

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 19.7462  | 5.2521     | 3.760   | 0.000765 | *** |
| wt          | -5.0480  | 0.4840     | -10.430 | 2.52e-11 | *** |
| qsec        | 0.9292   | 0.2650     | 3.506   | 0.001500 | **  |

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.596 on 29 degrees of freedom
```

```
Multiple R-squared:  0.8264, Adjusted R-squared:  0.8144
```

```
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12
```

# What's "messy" about a linear regression?

```
> summary(lmfit)
```

```
Call:
```

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max  
-4.3962 -2.1431 -0.2129  1.4915  5.7486
```

1. Extracting coefficients takes multiple steps:

```
data.frame(coef(summary(lmfit)))
```

```
Coefficients:
```

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 19.7462  | 5.2521     | 3.760   | 0.000765 | *** |
| wt          | -5.0480  | 0.4840     | -10.430 | 2.52e-11 | *** |
| qsec        | 0.9292   | 0.2650     | 3.506   | 0.001500 | **  |

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.596 on 29 degrees of freedom
```

```
Multiple R-squared:  0.8264, Adjusted R-squared:  0.8144
```

```
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12
```

# What's "messy" about a linear regression?

```
> summary(lmfit)
```

```
Call:
```

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-4.3962 -2.1431 -0.2129  1.4915  5.7486
```

```
Coefficients:
```

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 19.7462  | 5.2521     | 3.760   | 0.000765 | *** |
| wt          | -5.0480  | 0.4840     | -10.430 | 2.52e-11 | *** |
| qsec        | 0.9292   | 0.2650     | 3.506   | 0.001500 | **  |

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.596 on 29 degrees of freedom
```

```
Multiple R-squared:  0.8264, Adjusted R-squared:  0.8144
```

```
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12
```

2. Information in row names  
(can't combine models)

# What's "messy" about a linear regression?

```
> summary(lmfit)
```

```
Call:
```

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-4.3962 -2.1431 -0.2129  1.4915  5.7486
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  19.7462     5.2521    3.760 0.000765 ***
wt           -5.0480     0.4840   -10.430 2.52e-11 ***
qsec         0.9292     0.2650    3.506 0.001500 **
```

**3. Column names are inconvenient (access with `$"Pr(>|t|)"`, converts to `Pr...t...`)**

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.596 on 29 degrees of freedom
```

```
Multiple R-squared:  0.8264, Adjusted R-squared:  0.8144
```

```
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12
```

# What's "messy" about a linear regression?

```
> summary(lmfit)
```

```
Call:
```

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-4.3962 -2.1431 -0.2129  1.4915  5.7486
```

```
Coefficients:
```

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  19.7462     5.2521   3.760 0.000765 ***
wt           -5.0480     0.4840 -10.430 2.52e-11 ***
qsec         0.9292     0.2650   3.506 0.001500 **
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.596 on 29 degrees of freedom
```

```
Multiple R-squared:  0.8264, Adjusted R-squared:  0.8144
```

```
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12
```

**4. Information is computed in print method, not stored**

# What's "messy" about a linear regression?

```
> summary(lmfit)
```

Call:

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-4.3962 -2.1431 -0.2129  1.4915  5.7486
```

1. Extracting coefficients takes multiple steps:

```
data.frame(coef(summary(lmfit)))
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 19.7462  | 5.2521     | 3.760   | 0.000765 | *** |
| wt          | -5.0480  | 0.4840     | -10.430 | 2.52e-11 | *** |
| qsec        | 0.9292   | 0.2650     | 3.506   | 0.001500 | **  |

3. Column names are inconvenient (access with `$"Pr(>|t|)"`, converts to `Pr...t...`)

2. Information in row names (can't combine models)

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 2.596 on 29 degrees of freedom
```

```
Multiple R-squared:  0.8264, Adjusted R-squared:  0.8144
```

```
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12
```

4. Information is computed in `print` method, not stored

**These inconveniences aren't  
an exception, they're the rule**

broom's `tidy()` method  
does the work for you

```
> tidy(lmfit)
  term estimate std.error statistic p.value
1 (Intercept)  19.746    5.252      3.76 7.65e-04
2 wt          -5.048    0.484     -10.43 2.52e-11
3 qsec         0.929    0.265      3.51 1.50e-03
```

# broom's `tidy()` method does the work for you

**One function  
to call**

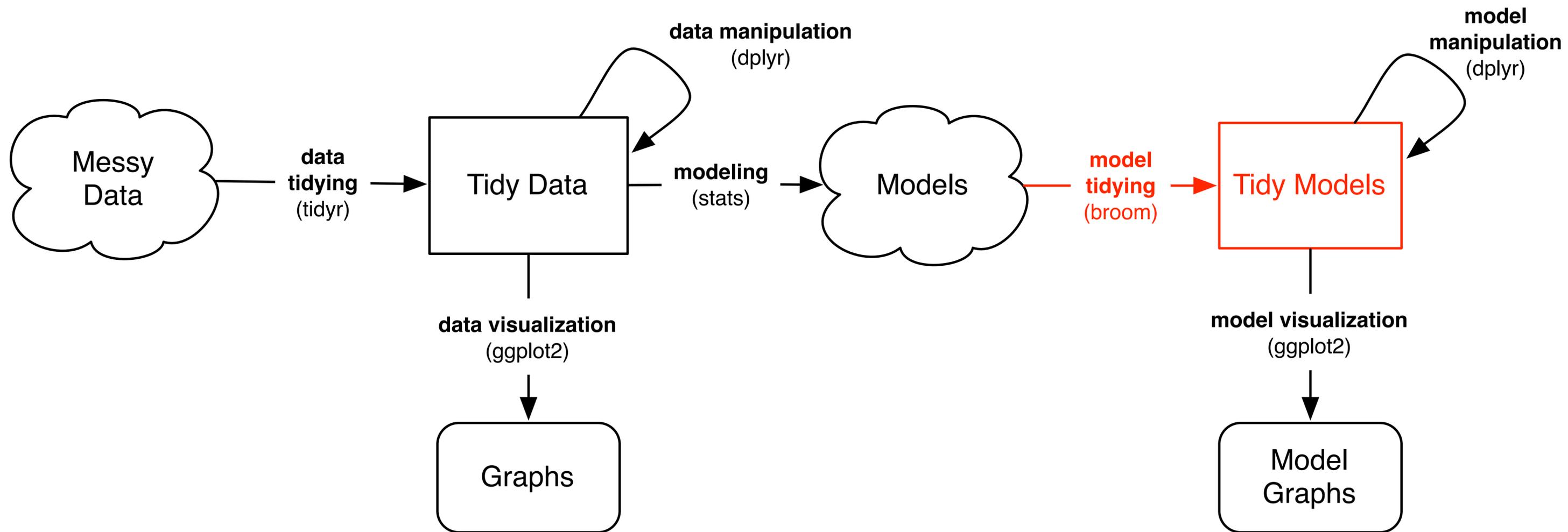
```
> tidy(lmfit)
```

|   | term        | estimate | std.error | statistic | p.value  |
|---|-------------|----------|-----------|-----------|----------|
| 1 | (Intercept) | 19.746   | 5.252     | 3.76      | 7.65e-04 |
| 2 | wt          | -5.048   | 0.484     | -10.43    | 2.52e-11 |
| 3 | qsec        | 0.929    | 0.265     | 3.51      | 1.50e-03 |

**Convenient  
column names**

**Information stored  
in columns, never  
row names**

broom takes model objects and  
turns them into tidy data frames  
that can be used with tidy tools



# broom's three methods

- broom defines tidying methods for extracting three kinds of statistics from an object:
  - **tidy()**: component-level statistics
  - **augment()**: observation-level statistics
  - **glance()**: model-level statistics

# Example: three levels of a linear regression

```
> summary(lmfit)
```

Call:

```
lm(formula = mpg ~ wt + qsec, data = mtcars)
```

**Observation Level:**  
fitted values, residuals  
`augment()`

Residuals:

| Min     | 1Q      | Median  | 3Q     | Max    |
|---------|---------|---------|--------|--------|
| -4.3962 | -2.1431 | -0.2129 | 1.4915 | 5.7486 |

**Component Level:**  
coefficients, p-values  
`tidy()`

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | 19.7462  | 5.2521     | 3.760   | 0.000765 | *** |
| wt          | -5.0480  | 0.4840     | -10.430 | 2.52e-11 | *** |
| qsec        | 0.9292   | 0.2650     | 3.506   | 0.001500 | **  |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

**Model Level:**  
 $R^2$ , F-statistic, deviance  
`glance()`

Residual standard error: 2.596 on 29 degrees of freedom  
Multiple R-squared: 0.8264, Adjusted R-squared: 0.8144  
F-statistic: 69.03 on 2 and 29 DF, p-value: 9.395e-12

# The `tidy()` method: component-level statistics

```
> tidy(lmfit)
  term estimate std.error statistic p.value
1 (Intercept) 19.746    5.252     3.76 7.65e-04
2 wt          -5.048    0.484    -10.43 2.52e-11
3 qsec         0.929    0.265     3.51 1.50e-03
```

← each row is a coefficient

# The `augment()` method: observation-level statistics

```
> augment(lmfit)
  .rownames mpg wt qsec .fitted .se.fit .resid .hat .sigma
1 Mazda RX4 21.0 2.62 16.5 21.82 0.683 -0.8151 0.0693 2.64
2 Mazda RX4 Wag 21.0 2.88 17.0 21.05 0.547 -0.0482 0.0444 2.64
3 Datsun 710 22.8 2.32 18.6 25.33 0.640 -2.5273 0.0607 2.60
4 Hornet 4 Drive 21.4 3.21 19.4 21.58 0.623 -0.1806 0.0576 2.64
5 Hornet Sportabout 18.7 3.44 17.0 18.20 0.512 0.5039 0.0389 2.64
6 Valiant 18.1 3.46 20.2 21.07 0.803 -2.9686 0.0957 2.58
7 Duster 360 14.3 3.57 15.8 16.44 0.701 -2.1434 0.0729 2.61
8 Merc 240D 24.4 3.19 20.0 22.23 0.730 2.1729 0.0791 2.61
9 Merc 230 22.8 3.15 22.9 25.12 1.410 -2.3237 0.2950 2.59
10 Merc 280 19.2 3.44 18.3 19.39 0.491 -0.1855 0.0358 2.64
11 Merc 280C 17.8 3.44 18.9 19.94 0.557 -2.1430 0.0460 2.61
12 Merc 450SE 16.4 4.07 17.4 15.37 0.615 1.0310 0.0561 2.63
13 Merc 450SL 17.3 3.73 17.6 17.27 0.520 0.0289 0.0402 2.64
14 Merc 450SLC 15.2 3.78 18.0 17.39 0.539 -2.1904 0.0431 2.61
15 Cadillac Fleetwood 10.4 5.25 18.0 9.95 1.092 0.4487 0.1768 2.64
16 Lincoln Continental 10.4 5.42 17.8 8.92 1.161 1.4757 0.2001 2.62
17 Chrysler Imperial 14.7 5.34 17.4 8.95 1.115 5.7486 0.1844 2.35
18 Fiat 128 32.4 2.20 19.5 26.73 0.751 5.6679 0.0836 2.39
19 Honda Civic 30.4 1.61 18.5 28.80 0.892 1.5975 0.1180 2.62
20 Toyota Corolla 33.9 1.83 19.9 28.97 0.909 4.9258 0.1226 2.45
```

each row is an  
observation from the  
original data

# The `augment()` method: observation-level statistics

note that new columns start with `.`

```
> augment(lmfit)
      .rownames  mpg   wt  qsec  .fitted  .se.fit  .resid  .hat  .sigma
1      Mazda RX4 21.0 2.62 16.5   21.82   0.683 -0.8151 0.0693 2.64
2      Mazda RX4 Wag 21.0 2.88 17.0   21.05   0.547 -0.0482 0.0444 2.64
3      Datsun 710 22.8 2.32 18.6   25.33   0.640 -2.5273 0.0607 2.60
4      Hornet 4 Drive 21.4 3.21 19.4   21.58   0.623 -0.1806 0.0576 2.64
5      Hornet Sportabout 18.7 3.44 17.0   18.20   0.512  0.5039 0.0389 2.64
6      Valiant 18.1 3.46 20.2   21.07   0.803 -2.9686 0.0957 2.58
7      Duster 360 14.3 3.57 15.8   16.44   0.701 -2.1434 0.0729 2.61
8      Merc 240D 24.4 3.19 20.0   22.23   0.730  2.1729 0.0791 2.61
9      Merc 230 22.8 3.15 22.9   25.12   1.410 -2.3237 0.2950 2.59
10     Merc 280 19.2 3.44 18.3   19.39   0.491 -0.1855 0.0358 2.64
11     Merc 280C 17.8 3.44 18.9   19.94   0.557 -2.1430 0.0460 2.61
12     Merc 450SE 16.4 4.07 17.4   15.37   0.615  1.0310 0.0561 2.63
13     Merc 450SL 17.3 3.73 17.6   17.27   0.520  0.0289 0.0402 2.64
14     Merc 450SLC 15.2 3.78 18.0   17.39   0.539 -2.1904 0.0431 2.61
15     Cadillac Fleetwood 10.4 5.25 18.0    9.95   1.092  0.4487 0.1768 2.64
16     Lincoln Continental 10.4 5.42 17.8    8.92   1.161  1.4757 0.2001 2.62
17     Chrysler Imperial 14.7 5.34 17.4    8.95   1.115  5.7486 0.1844 2.35
18     Fiat 128 32.4 2.20 19.5   26.73   0.751  5.6679 0.0836 2.39
19     Honda Civic 30.4 1.61 18.5   28.80   0.892  1.5975 0.1180 2.62
20     Toyota Corolla 33.9 1.83 19.9   28.97   0.909  4.9258 0.1226 2.45
```

each row is an  
observation from the  
original data

# The `glance()` method: model-level statistics

```
> glance(lmfit)
  r.squared adj.r.squared sigma statistic  p.value  df logLik AIC  BIC deviance
1    0.826      0.814    2.6         69 9.39e-12  3  -74.4 157 163    195 ← one row for the model
```

broom works across many  
kinds of model objects

# Nonlinear least squares: **before**

```
> n <- nls(mpg ~ k * e ^ wt, data = mtcars, start = list(k = 1, e = 2))  
> summary(n)
```

Formula: mpg ~ k \* e^wt

Parameters:

|   | Estimate | Std. Error | t value | Pr(> t ) |     |
|---|----------|------------|---------|----------|-----|
| k | 49.6597  | 3.7888     | 13.1    | 6e-14    | *** |
| e | 0.7456   | 0.0199     | 37.5    | <2e-16   | *** |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.67 on 30 degrees of freedom

Number of iterations to convergence: 10

Achieved convergence tolerance: 2.04e-06

# Nonlinear least squares: after

```
> tidy(n)
  term estimate std.error statistic  p.value
1    k   49.660    3.7888     13.1 5.96e-14 ← each row is one estimated parameter
2    e    0.746    0.0199     37.5 8.86e-27

> augment(n)
  mpg  wt  .fitted .resid
1 21.0 2.62   23.0 -2.012 ← each row is an observation from the original data
2 21.0 2.88   21.4 -0.352
3 22.8 2.32   25.1 -2.331
4 21.4 3.21   19.3  2.076
5 18.7 3.44   18.1  0.611
6 18.1 3.46   18.0  0.117
...

> glance(n)
  sigma isConv  finTol logLik AIC BIC deviance df.residual
1  2.67  TRUE 2.04e-06  -75.8 158 162      214      30 ← one row for the model
```



# K-means clustering: after

```
> tidy(k)
  x1      x2      x3      x4 size withinss cluster
1 5.901613 2.748387 4.393548 1.433871 62 39.82097 1
2 5.006000 3.428000 1.462000 0.246000 50 15.15100 2
3 6.850000 3.073684 5.742105 2.071053 38 23.87947 3
```

← each row is one cluster

```
> head(augment(k, m))
  Sepal.Length Sepal.Width Petal.Length Petal.Width .cluster
1           5.1           3.5           1.4           0.2           2
2           4.9           3.0           1.4           0.2           2
3           4.7           3.2           1.3           0.2           2
4           4.6           3.1           1.5           0.2           2
5           5.0           3.6           1.4           0.2           2
6           5.4           3.9           1.7           0.4           2
```

← each row is one assignment

```
> glance(k)
  totss tot.withinss betweenss iter
1 681.3706      78.85144  602.5192  2
```

← one row describing the entire clustering operation

# And many others...

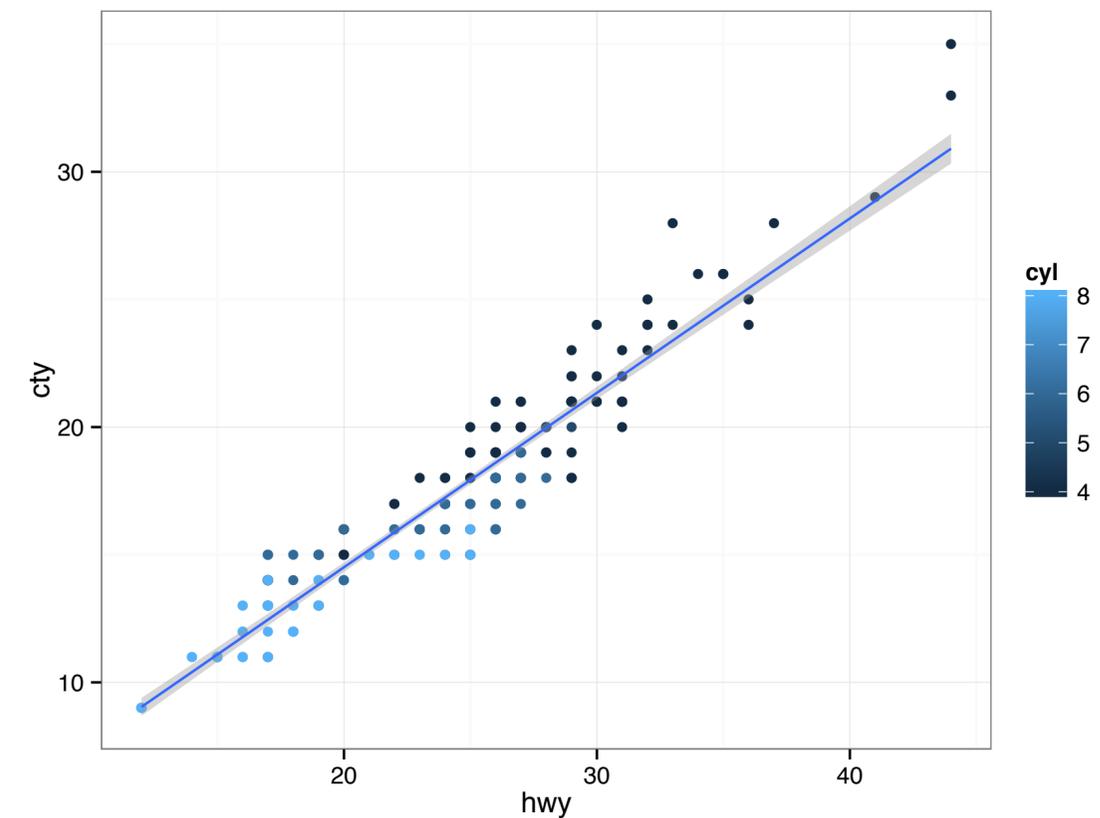
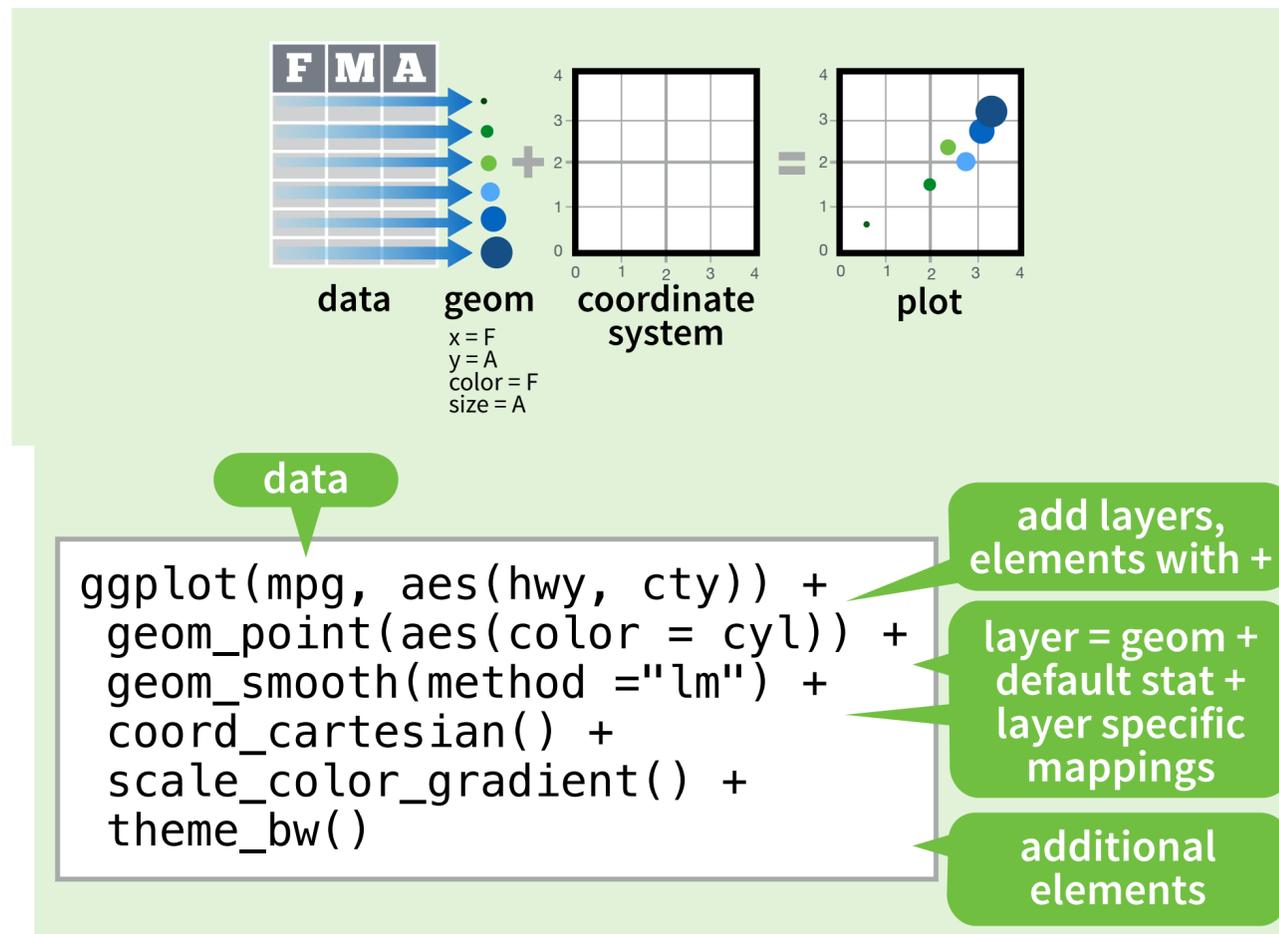
| Class       | tidy | glance | augment |
|-------------|------|--------|---------|
| aareg       | x    | x      |         |
| anova       | x    |        |         |
| aov         | x    |        |         |
| aovlist     | x    |        |         |
| Arima       | x    | x      |         |
| biglm       | x    | x      |         |
| binDesign   | x    | x      |         |
| binWidth    | x    |        |         |
| boot        | x    |        |         |
| btergm      | x    |        |         |
| cch         | x    | x      |         |
| cid         | x    |        |         |
| coefest     | x    |        |         |
| confint.glm | x    |        |         |
| coxph       | x    | x      | x       |
| cv.glmnet   | x    | x      |         |
| data.frame  | x    | x      | x       |
| default     | x    | x      | x       |
| density     | x    |        |         |
| ergm        | x    | x      |         |
| felm        | x    | x      | x       |
| fitdistr    | x    | x      |         |
| ftable      | x    |        |         |
| gam         | x    | x      |         |
| gamiss      | x    |        |         |
| geeglm      | x    |        |         |
| glht        | x    |        |         |

| Class          | tidy | glance | augment |
|----------------|------|--------|---------|
| glmnet         | x    | x      |         |
| gmm            | x    | x      |         |
| htest          | x    | x      |         |
| kappa          | x    |        |         |
| kmeans         | x    | x      | x       |
| Line           | x    |        |         |
| Lines          | x    |        |         |
| list           | x    | x      |         |
| lm             | x    | x      | x       |
| lme            | x    | x      | x       |
| manova         | x    |        |         |
| map            | x    |        |         |
| matrix         | x    | x      |         |
| merMod         | x    | x      | x       |
| mle2           | x    |        |         |
| multinom       | x    | x      |         |
| nlrq           | x    | x      | x       |
| nls            | x    | x      | x       |
| NULL           | x    | x      | x       |
| pairwise.htest | x    |        |         |
| plm            | x    | x      | x       |
| Polygon        | x    |        |         |
| Polygons       | x    |        |         |
| power.htest    | x    |        |         |
| pyears         | x    | x      |         |
| rcorr          | x    |        |         |
| ridgelm        | x    | x      |         |

| Class                    | tidy | glance | augment |
|--------------------------|------|--------|---------|
| rjags                    | x    |        |         |
| roc                      | x    |        |         |
| rowwise_df               | x    | x      | x       |
| rq                       | x    | x      | x       |
| rqs                      | x    | x      | x       |
| SpatialLinesDataFrame    | x    |        |         |
| SpatialPolygons          | x    |        |         |
| SpatialPolygonsDataFrame | x    |        |         |
| spec                     | x    |        |         |
| stanfit                  | x    |        |         |
| summary.glm              | x    |        |         |
| summaryDefault           | x    | x      |         |
| survexp                  | x    | x      |         |
| survfit                  | x    | x      |         |
| survreg                  | x    | x      | x       |
| table                    | x    |        |         |
| tbl_df                   | x    | x      | x       |
| ts                       | x    |        |         |
| TukeyHSD                 | x    |        |         |
| zoo                      | x    |        |         |

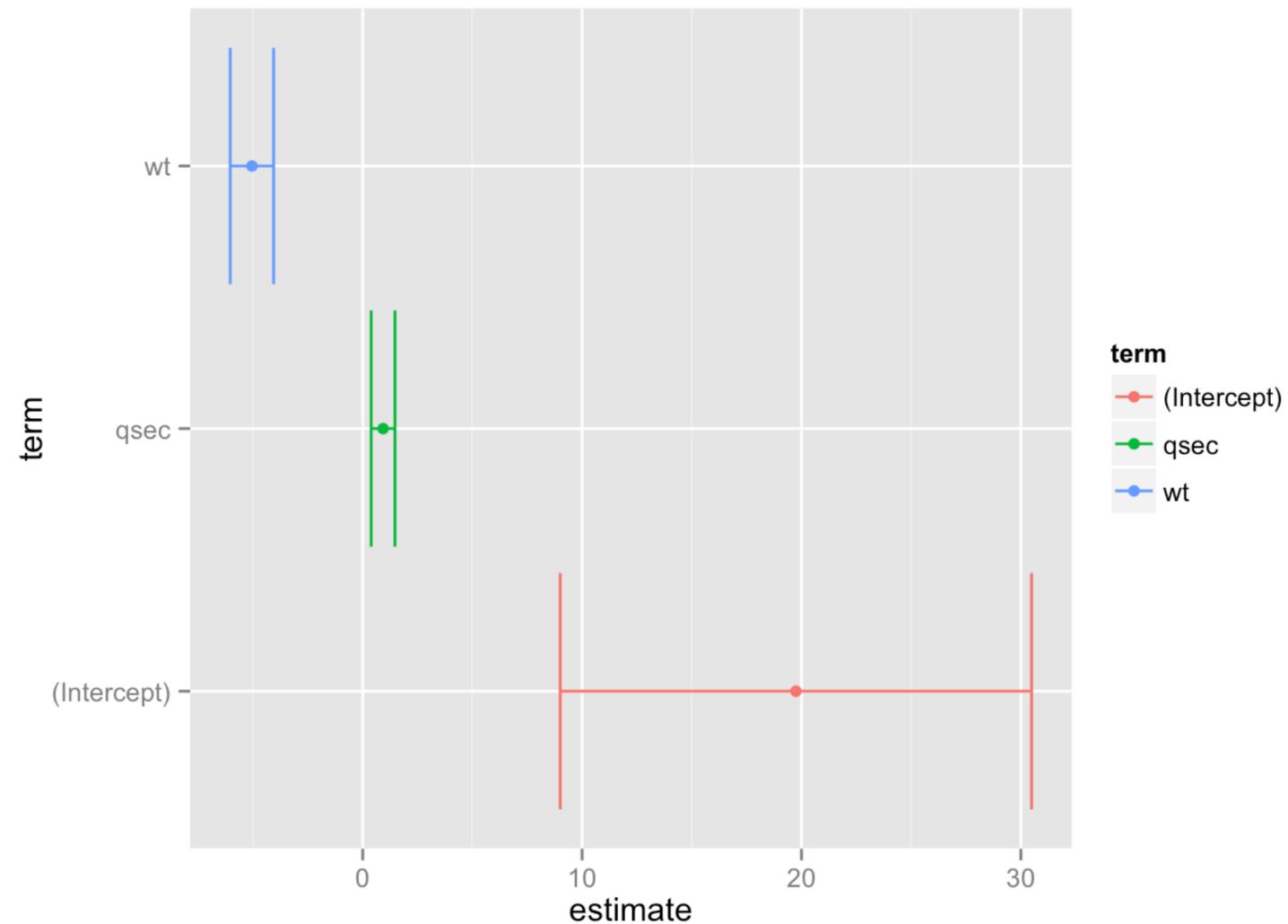
Why are tidy models  
useful?

# ggplot2 can visualize tidy data



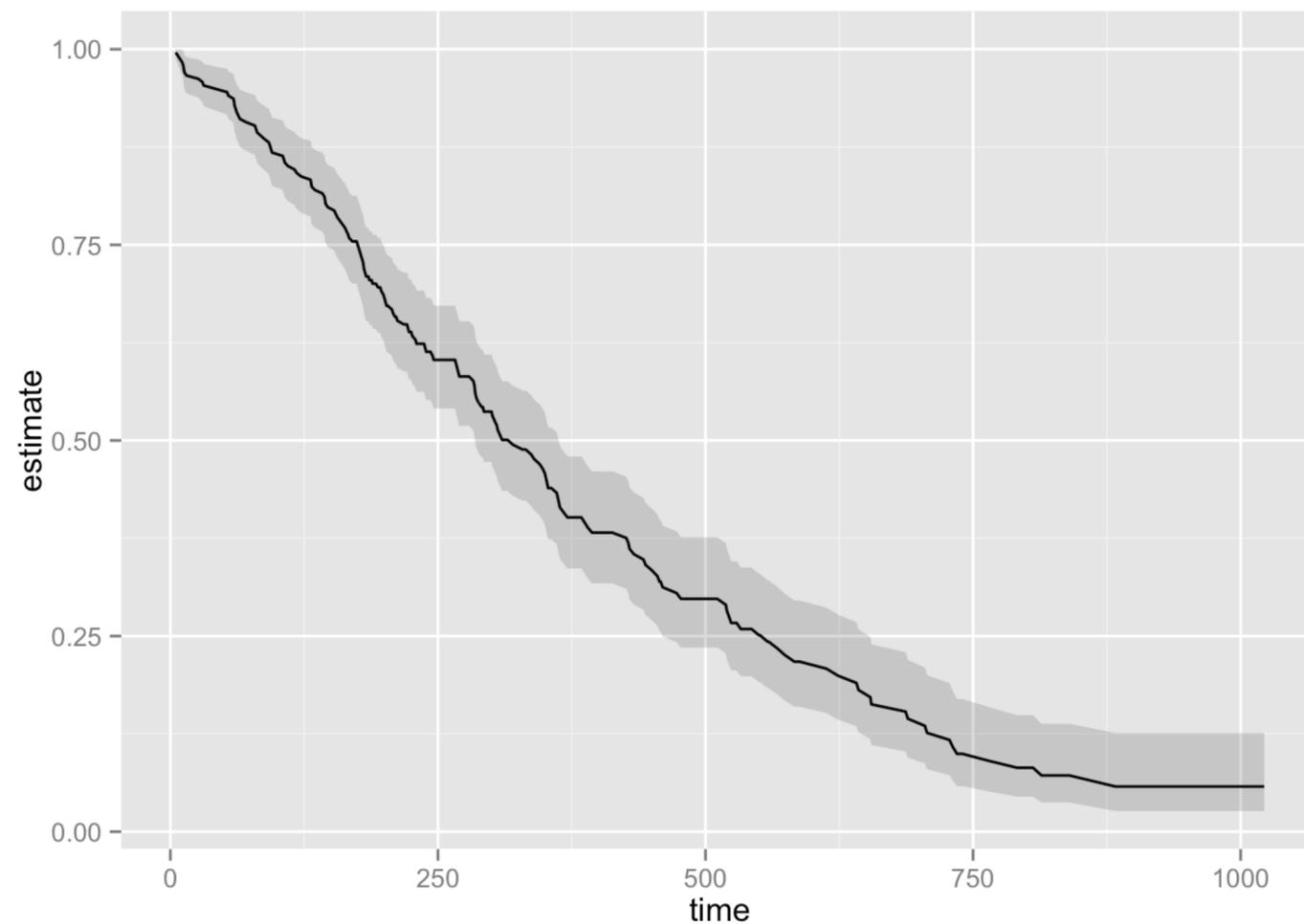
# Example: coefficient plot

```
td <- tidy(lmfit, conf.int = TRUE)
ggplot(td, aes(estimate, term, color = term)) +
  geom_point() +
  geom_errorbarh(aes(xmin = conf.low, xmax = conf.high))
```



# Example: survival curves

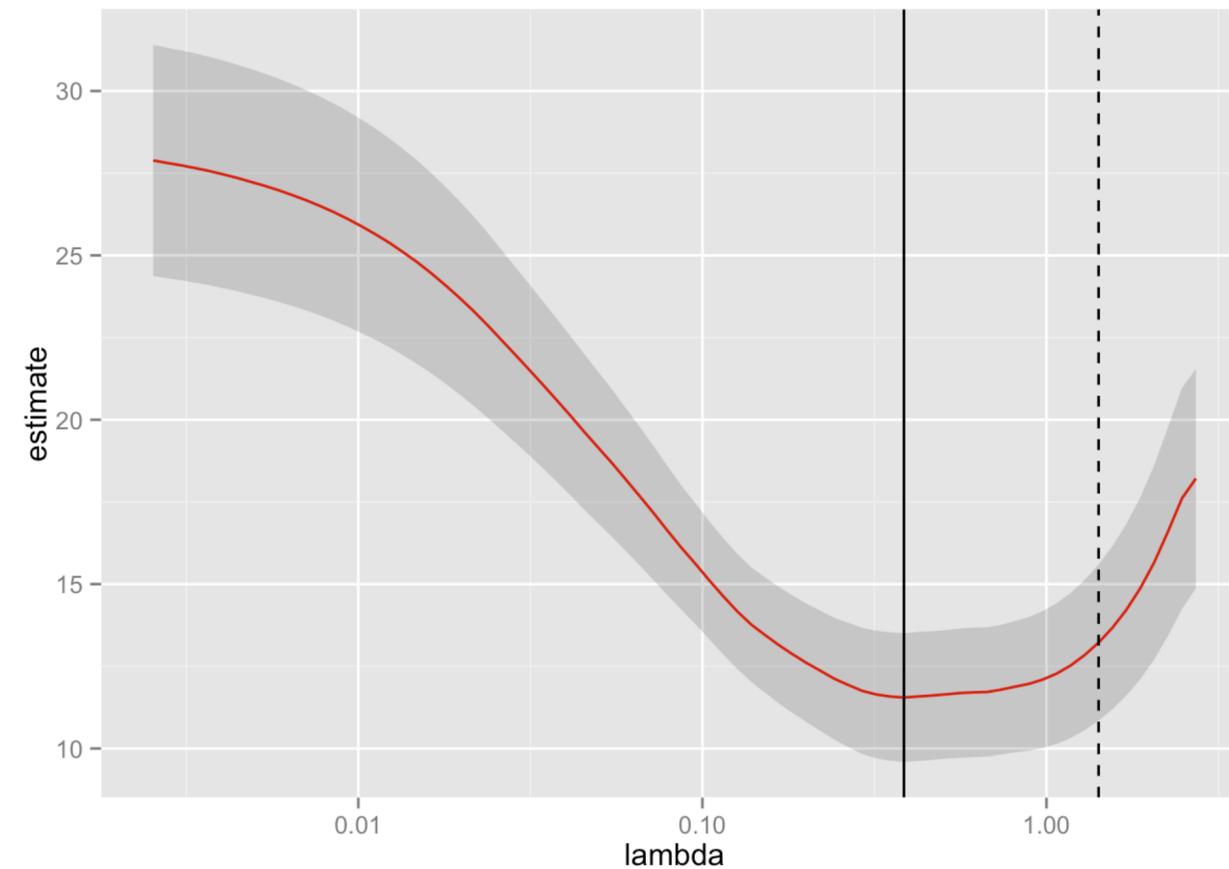
```
library(survival)
surv_fit <- survfit(coxph(Surv(time, status) ~ age + sex, lung))
td <- tidy(surv_fit)
ggplot(td, aes(time, estimate)) + geom_line() +
  geom_ribbon(aes(ymin = conf.low, ymax = conf.high), alpha = .2)
```



# Example: LASSO regression

```
tidied_cv <- tidy(glmnet_fit)
glance_cv <- glance(glmnet_fit)

ggplot(tidied_cv, aes(lambda, estimate)) + geom_line(color = "red") +
  geom_ribbon(aes(ymin = conf.low, ymax = conf.high), alpha = .2) +
  scale_x_log10() +
  geom_vline(xintercept = glance_cv$lambda.min) +
  geom_vline(xintercept = glance_cv$lambda.1se, lty = 2)
```



# Multiple Models

# Tidy models can be combined and compared

|              |
|--------------|
| <b>Model</b> |

|   | term        | estimate | std.error | statistic | p.value  |
|---|-------------|----------|-----------|-----------|----------|
| 1 | (Intercept) | 19.746   | 5.252     | 3.76      | 7.65e-04 |
| 2 | wt          | -5.048   | 0.484     | -10.43    | 2.52e-11 |
| 3 | qsec        | 0.929    | 0.265     | 3.51      | 1.50e-03 |

- different parameters
- different methods
- bootstrap replicates
- subgroup models (within each country, gene...)
- ensemble voting

# Tidy models can be combined and compared

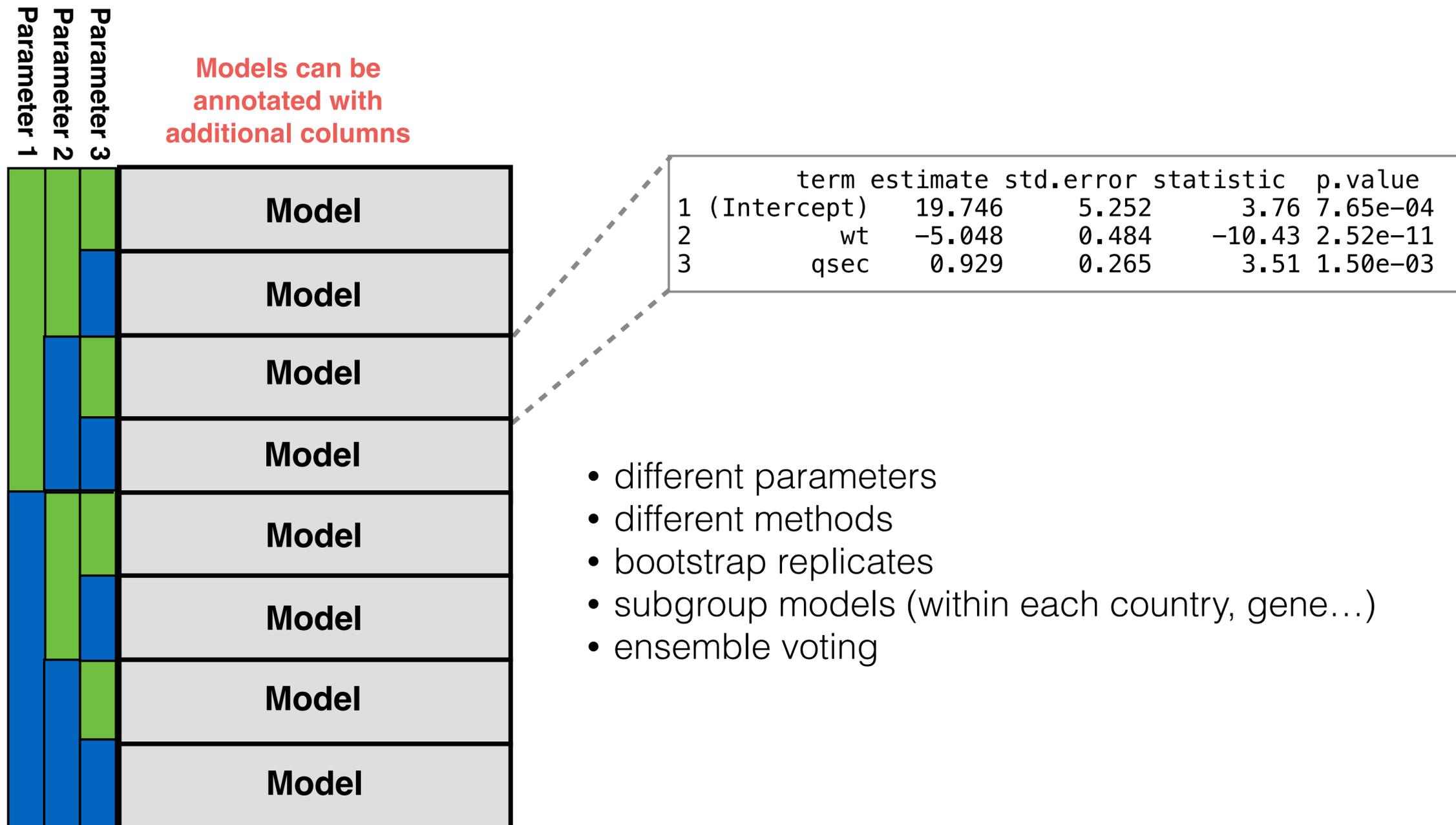
Parameter

Models can be annotated with additional columns

| Model | term        | estimate | std.error | statistic | p.value  |
|-------|-------------|----------|-----------|-----------|----------|
| 1     | (Intercept) | 19.746   | 5.252     | 3.76      | 7.65e-04 |
| 2     | wt          | -5.048   | 0.484     | -10.43    | 2.52e-11 |
| 3     | qsec        | 0.929    | 0.265     | 3.51      | 1.50e-03 |

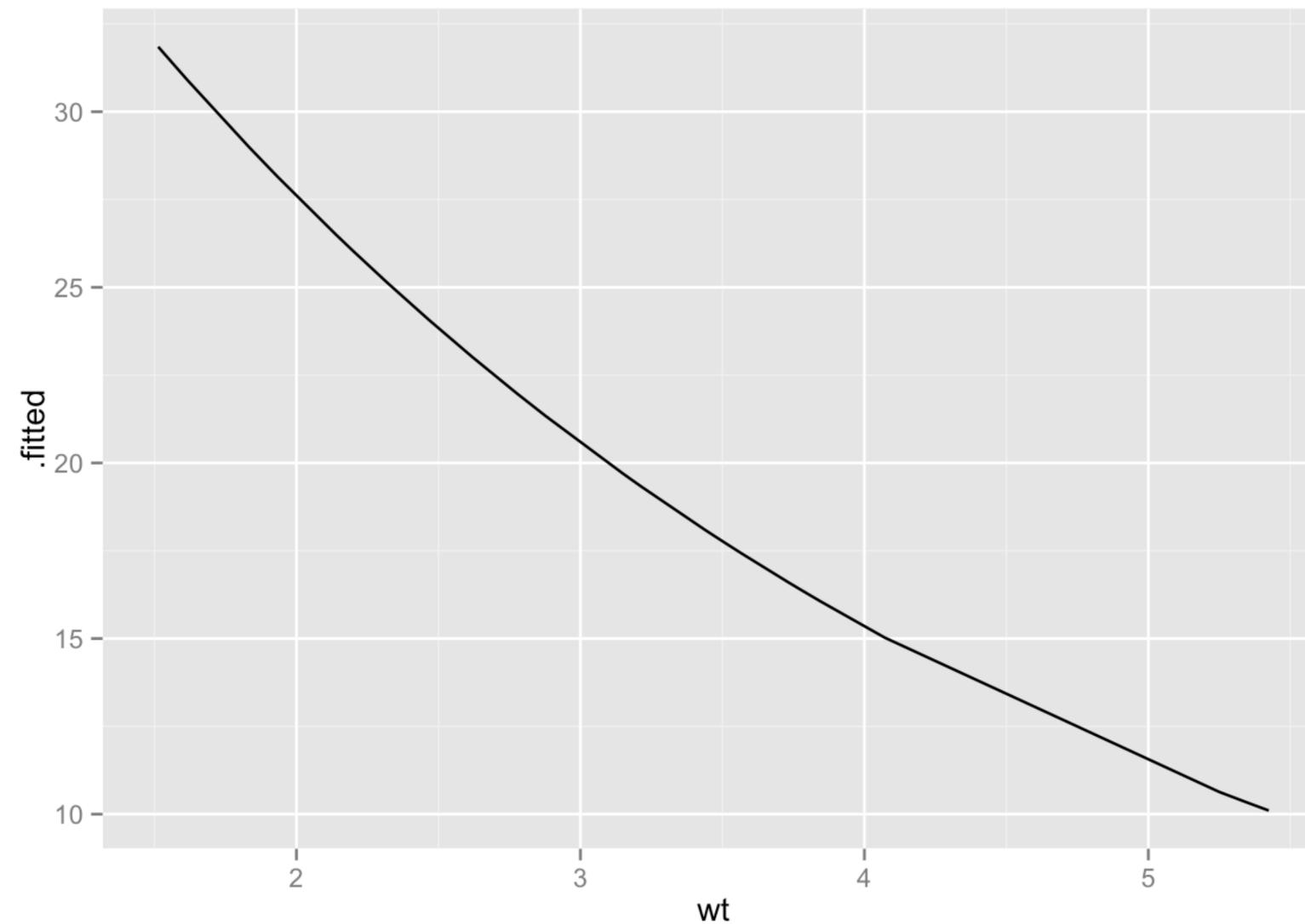
- different parameters
- different methods
- bootstrap replicates
- subgroup models (within each country, gene...)
- ensemble voting

# Tidy models can be combined and compared



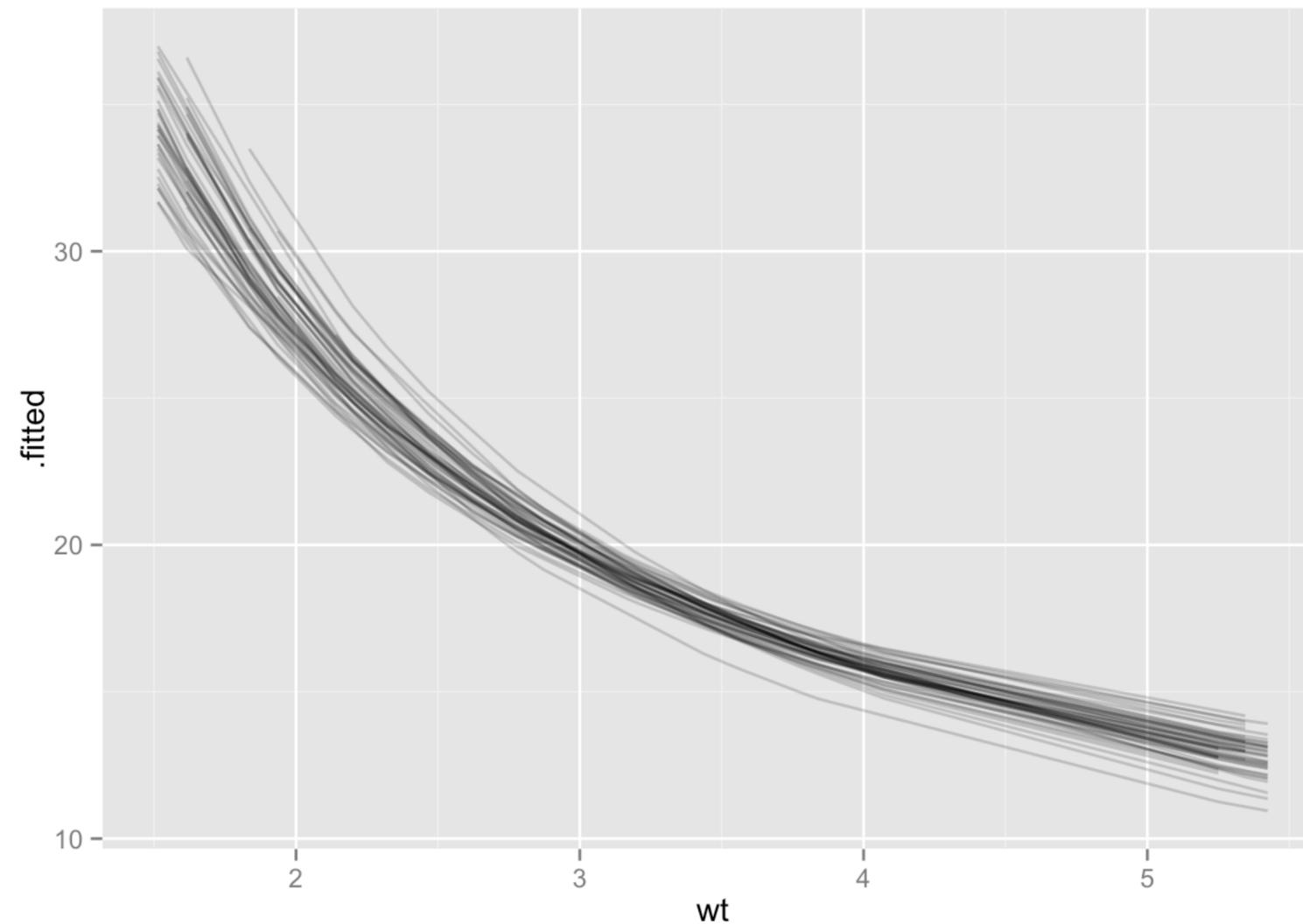
# If you can plot one nonlinear least squares fit...

```
augmented <- augment(nlsfit)  
ggplot(augmented, aes(wt, .fitted)) + geom_line()
```



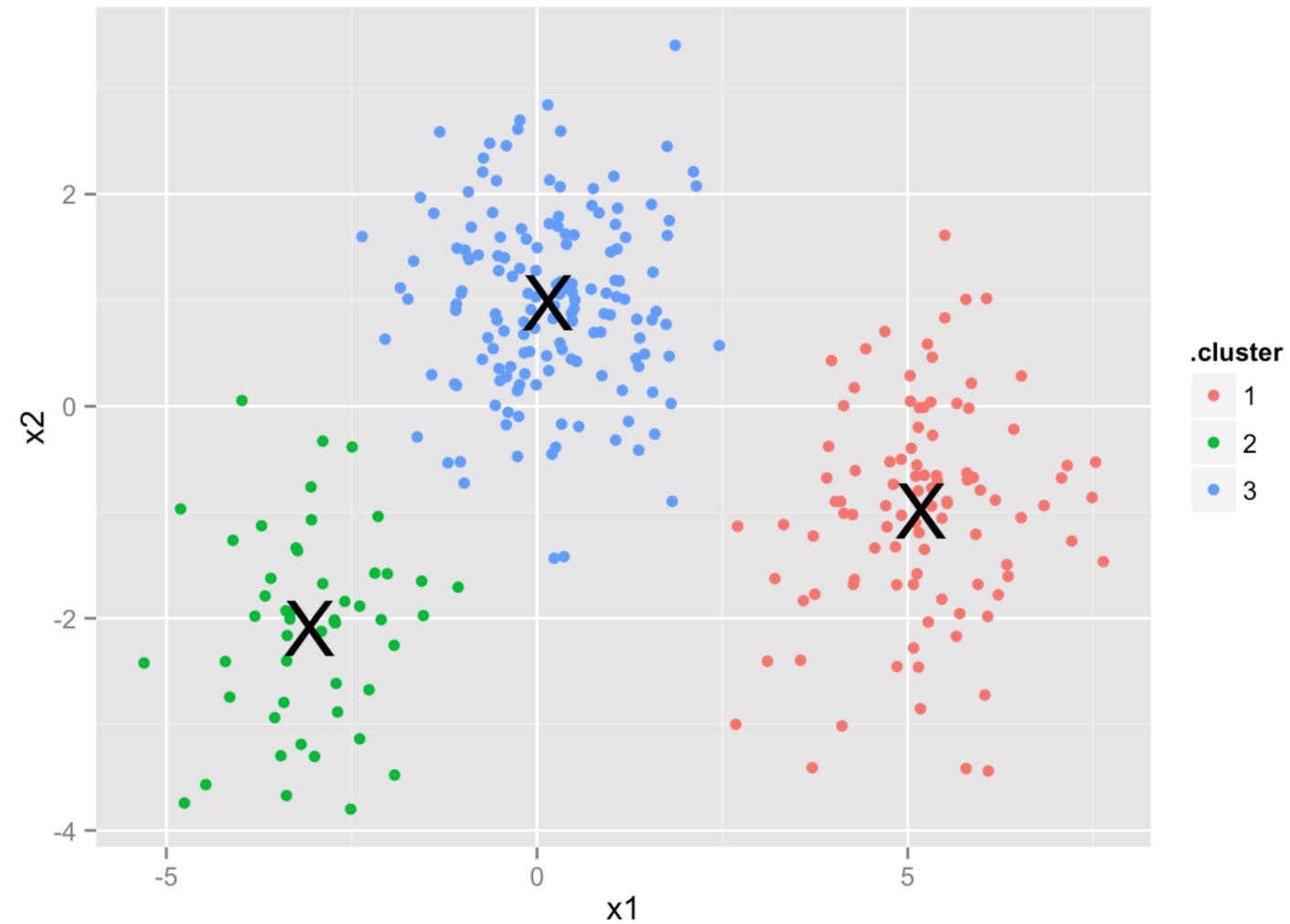
...you can plot 50 bootstrap replicates of it

```
ggplot(combined_augmented, aes(wt, .fitted, group = replicate)) +  
  geom_line(alpha = .2)
```



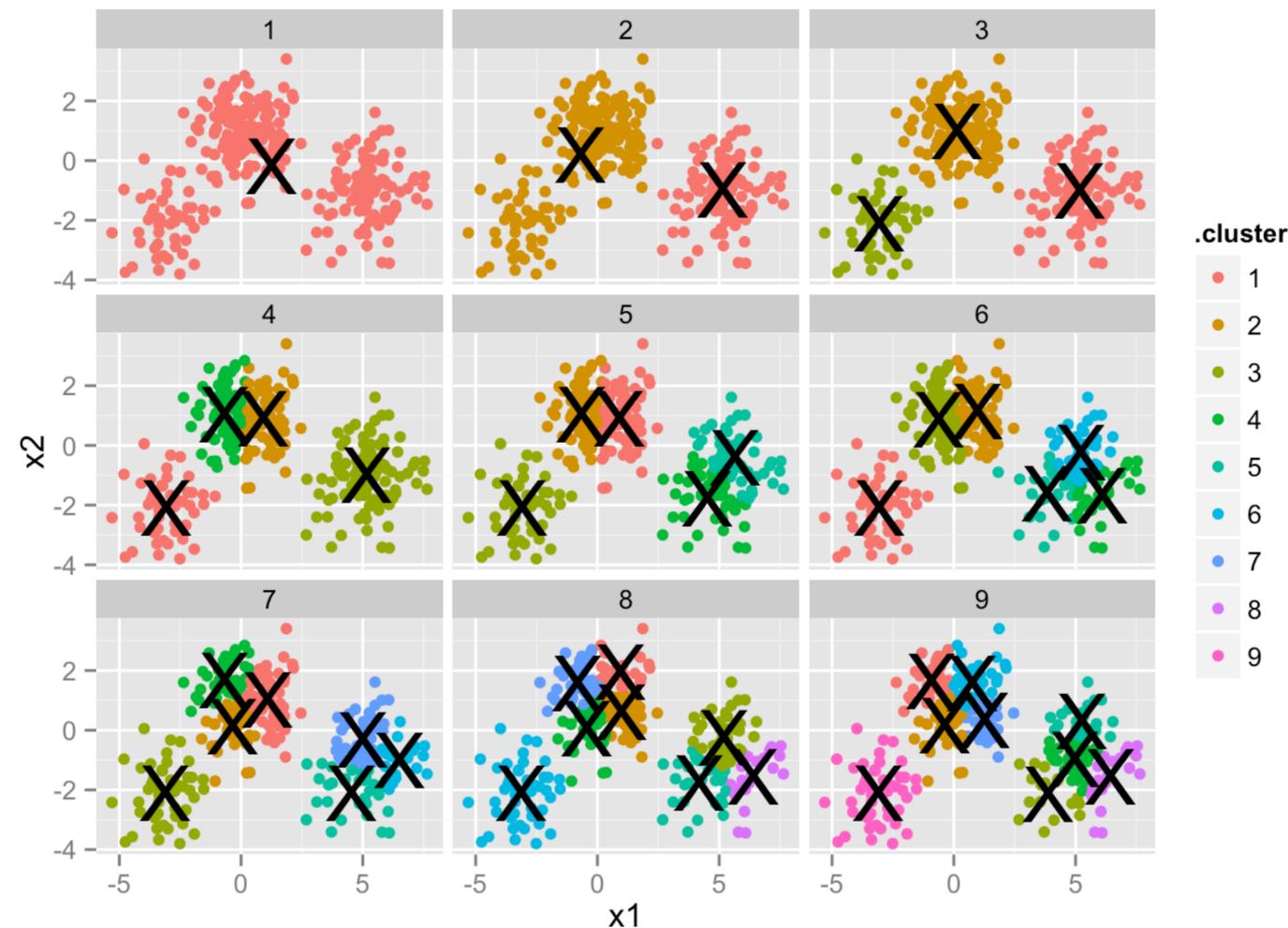
# If you can plot one instance of k-means clustering...

```
ggplot(assignments, aes(x1, x2)) +  
  geom_point(aes(color = .cluster)) +  
  geom_point(data = clusters, size = 10, shape = "x")
```

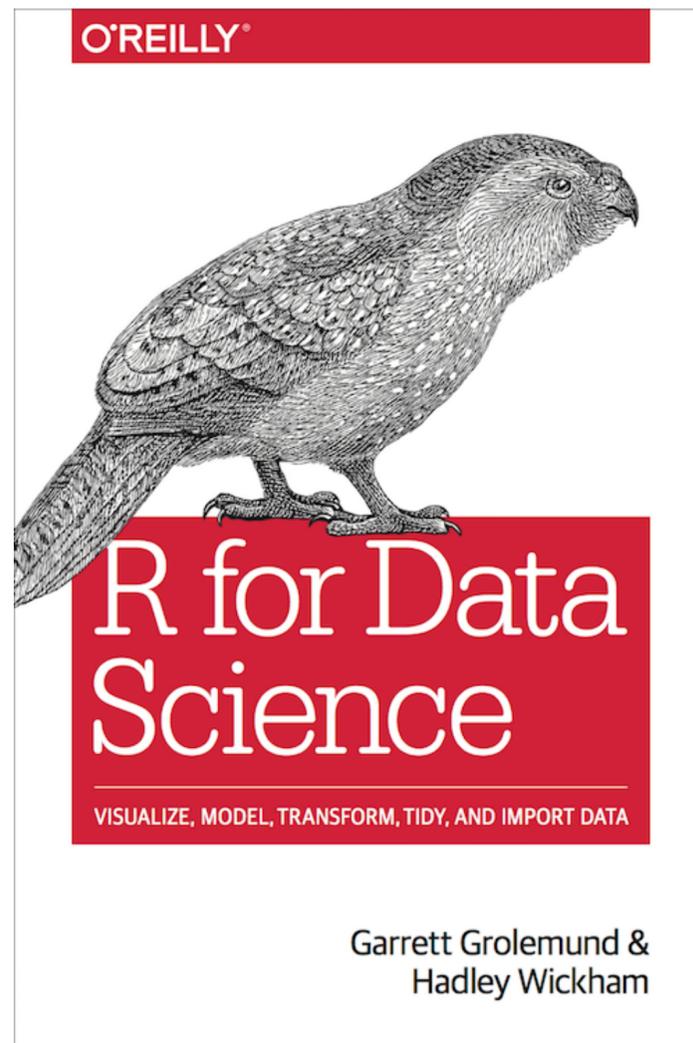


...you can plot it for many values of **k**

```
ggplot(combined_assignments, aes(x1, x2)) +  
  geom_point(aes(color = .cluster)) +  
  geom_point(data = combined_clusters, size = 10, shape = "x") +  
  facet_wrap(~ k)
```



# Learn more: many models



☰ 🔍 A ✎ R for Data Science

## 21 Many models

In this chapter you're going to learn three powerful ideas that help you to work with large numbers of models with ease:

1. Using many simple models to better understand complex datasets.
2. Using list-columns to store arbitrary data structures in a data frame. For example, this will allow you to have a column that contains linear models.
3. Using the **broom** package, by David Robinson, to turn models into tidy data. This is a powerful technique for working with large numbers of models because once you have tidy data, you can apply all of the techniques that you've learned about in earlier in the book.

<http://r4ds.had.co.nz/>

# Learn more: vignettes

[Introduction to broom](#)

[broom and dplyr](#)

[kmeans with dplyr+broom](#)

[Tidy bootstrapping with dplyr+broom](#)

# Learn more: manuscript

<http://arxiv.org/pdf/1412.3565v2.pdf>

## **broom**: An R Package for Converting Statistical Analysis Objects Into Tidy Data Frames

David Robinson

### **Abstract**

The concept of "tidy data" offers a powerful framework for structuring data to ease manipulation, modeling and visualization. However, most R functions, both those built-in and those found in third-party packages, produce output that is not tidy, and that is therefore difficult to reshape, recombine, and otherwise manipulate. Here I introduce the **broom** package, which turns the output of model objects into tidy data frames that are suited to further analysis, manipulation, and visualization with input-tidy tools. **broom** defines the `tidy`, `augment`, and `glance` generics, which arrange a model into three levels of tidy output respectively: the component level, the observation level, and the model level. I provide examples to demonstrate how these generics work with tidy tools to allow analysis and modeling of data that is divided into subsets, to recombine results from bootstrap replicates, and to perform simulations that investigate the effect of varying input parameters.

# Contribute: GitHub

 **dgrtwo** / **broom**

 Unwatch 37  Unstar 406  Fork 71

[Code](#) [Issues 41](#) [Pull requests 1](#) [Wiki](#) [Pulse](#) [Graphs](#) [Settings](#)

Convert statistical analysis objects from R into tidy format — Edit

 **332** commits  **2** branches  **9** releases  **25** contributors

Branch: **master** [New pull request](#) [Create new file](#) [Upload files](#) [Find file](#) [Clone or download](#)

 **dgrtwo** Version bump to prepare for CRAN Latest commit 35b7b46 3 days ago

|   |  |              |
|---|--|--------------|
|  <a href="#">R</a>             | Moved acf tidiers to stats_tidiers.                                      | 4 days ago   |
|  <a href="#">inst/extdata</a>  | Various edits to MCMC tidiers; mostly style changes. Added 8schools.s... | 7 months ago |
|  <a href="#">man-roxygen</a>   | Overhaul of how augmenting works across many objects. In particular t... | 2 years ago  |
|  <a href="#">man</a>           | Moved acf tidiers to stats_tidiers.                                      | 4 days ago   |
|  <a href="#">tests</a>         | Fixed to be compatible with dplyr 0.5                                    | 4 days ago   |
|  <a href="#">vignettes</a>     | update bootstrap vignette  | 4 months ago |
|  <a href="#">.Rbuildignore</a> | Added codecov.io   | 4 days ago   |
|  <a href="#">.gitignore</a>    | Update cran comments.  | 2 years ago  |
|  <a href="#">.travis.yml</a>   | Added codecov.io   | 4 days ago   |
|  <a href="#">CONDUCT.md</a>    | Added Code of Conduct  | 2 months ago |

# Thank you!

- broom package/paper
  - Matthieu Gomez
  - Boris Demeshev
  - Dieter Meine
  - Benjamin Nutter
  - Luke Johnston
  - Ben Bolker
  - Francois Briatte
  - Bob Muenchen
  - Hadley Wickham