# LINEAR MODELS ARE GREAT
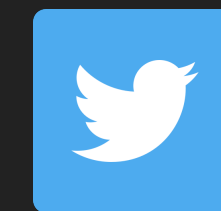
ALEX GOLD
SOLUTIONS ENGINEER
RSTUDIO

@alexkgold

Slides at: https://github.com/akgold/dsdc_linear_models

21.0
21.0
22.8
21.4
18.7
18.1
14.3
24.4
22.8
19.2
17.8
16.4
17.3
15.2
10.4
10.4
14.7
32.4
30.4
33.9
21.5
15.5
15.2
13.3
19.2
27.3
26.0
30.4
15.8
19.7
15.0
21.4

$$= \beta_0 + \beta_{cyl}$$

Just a number

6
6
4
6
8
6
8
4
4
6
6
8
8
8
8
8
8
4
4
4
4
8
8
8
8
4
4
4
8
6
8
4

$$+\beta_{hp}$$

160.0
160.0
108.0
258.0
360.0
225.0
360.0
146.7
140.8
167.6
167.6
275.8
275.8
275.8
472.0
460.0
440.0
78.7
75.7
71.1
120.1
318.0
304.0
350.0
400.0
79.0
120.3
95.1
351.0
145.0
301.0
121.0

+ some prediction error

$$Y = \beta X + \epsilon$$

# ALEX <3 LINEAR MODELS

```
> lm(mpg ~ cyl + disp + hp, data = mtcars)

Call:
lm(formula = mpg ~ cyl + disp + hp, data = mtcars)

Coefficients:
(Intercept)          cyl          disp            hp
   34.18492     -1.22742      -0.01884      -0.01468
```
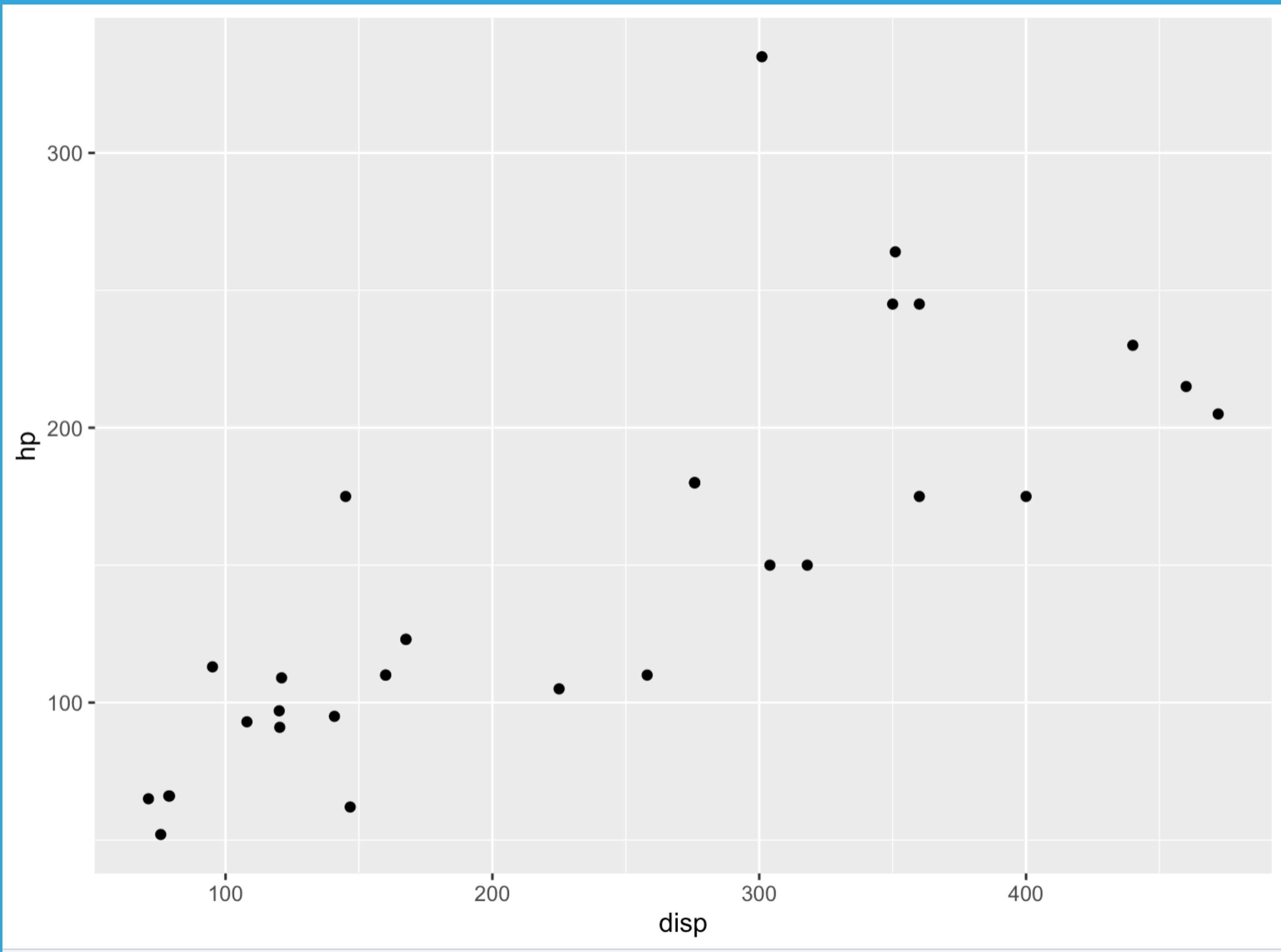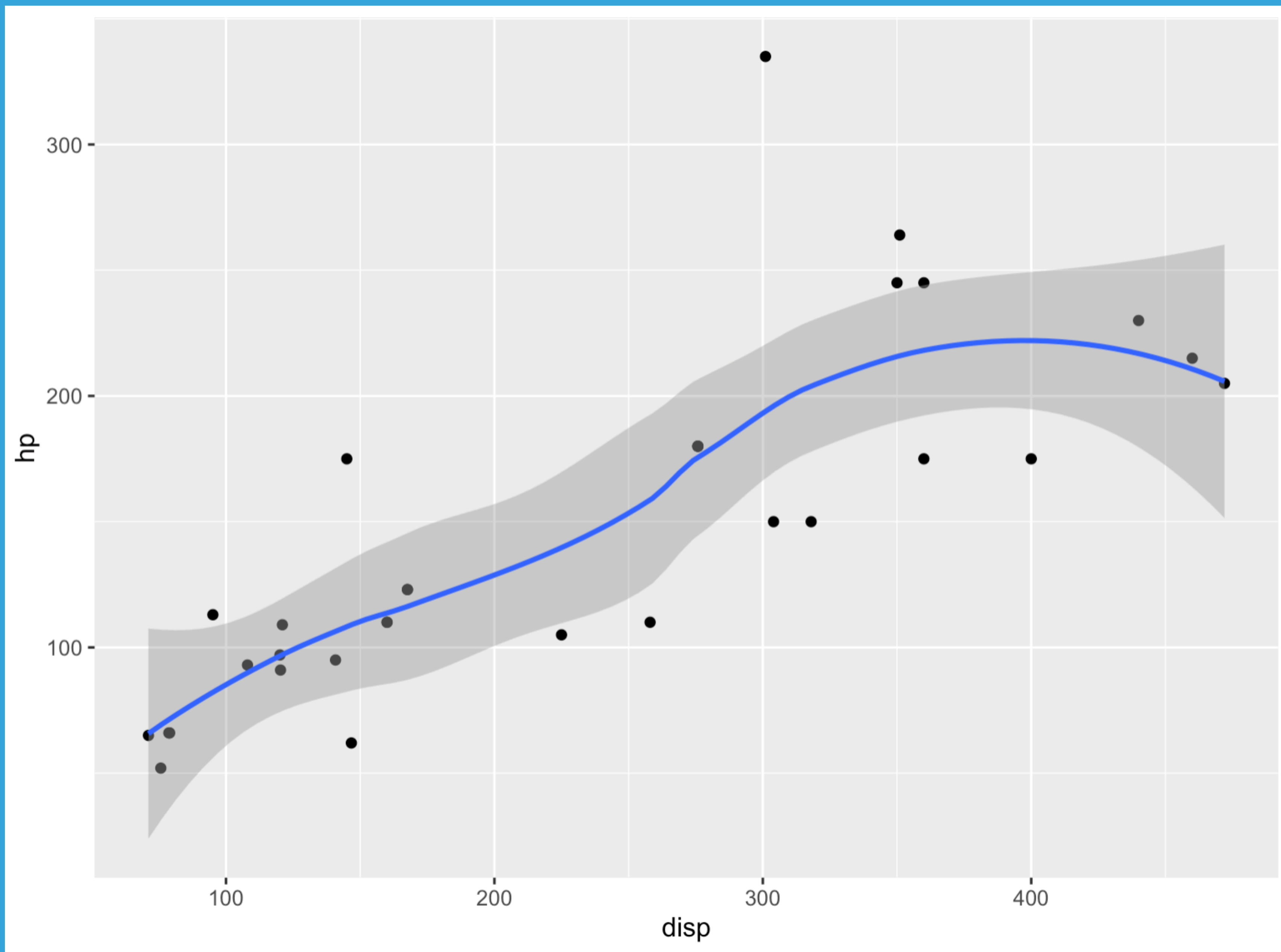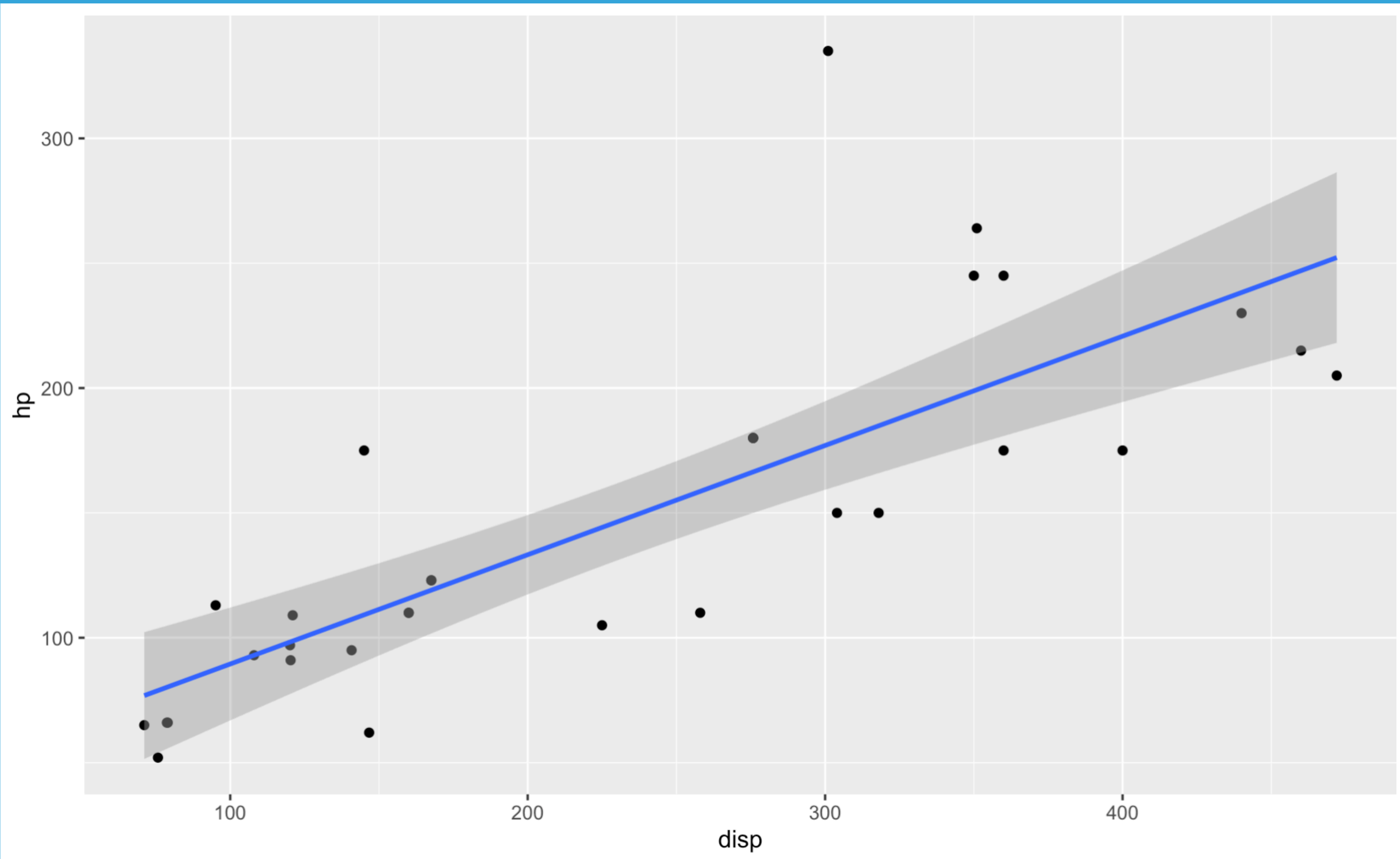
$$Y = \beta X + \epsilon$$

$$Y = \sum_k \beta_k f_k(x_i) + \epsilon$$

```
ggplot(mtcars, aes(x = disp, y = hp)) + geom_point() + stat_smooth(method = "loess")
```
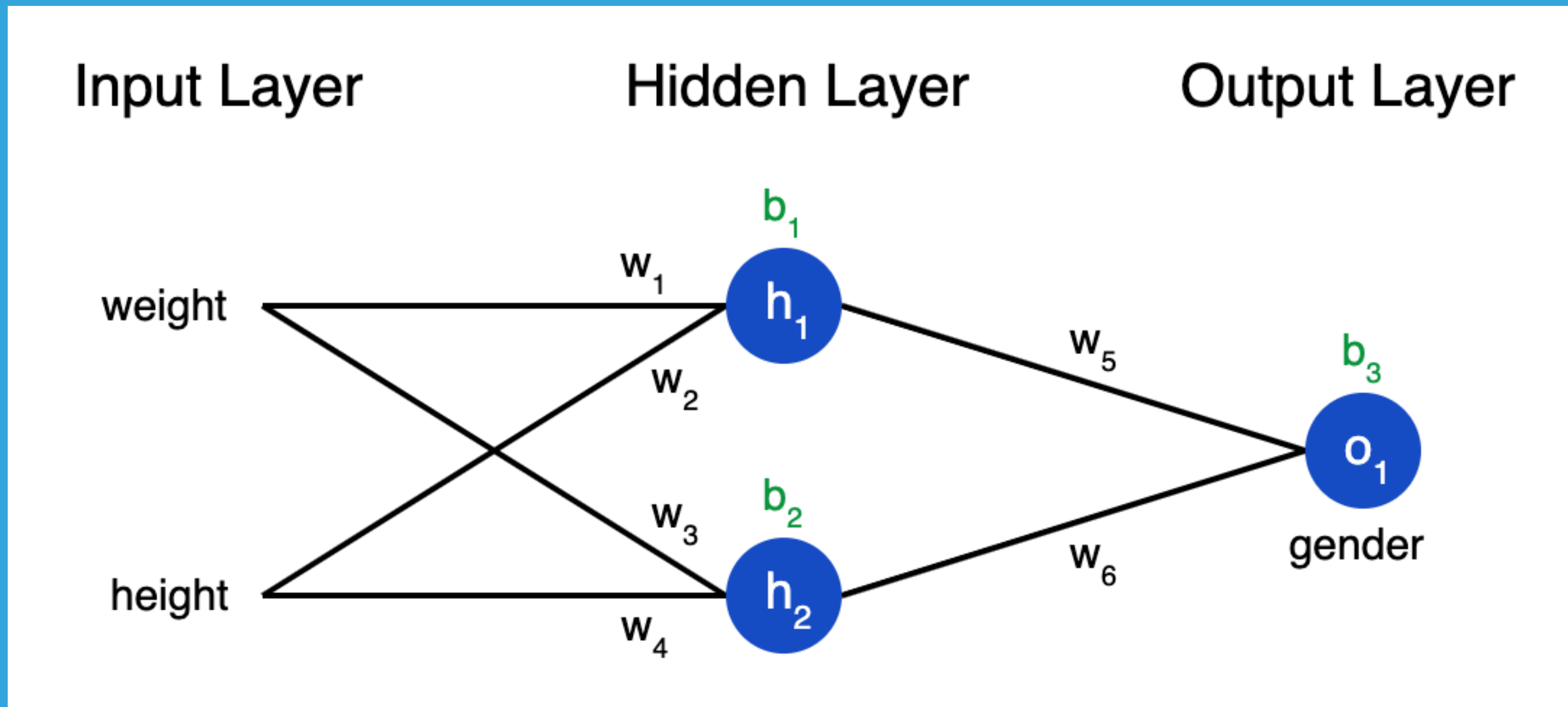
```
ggplot(mtcars, aes(x = disp, y = hp)) + geom_point() + stat_smooth(method = "lm")
```

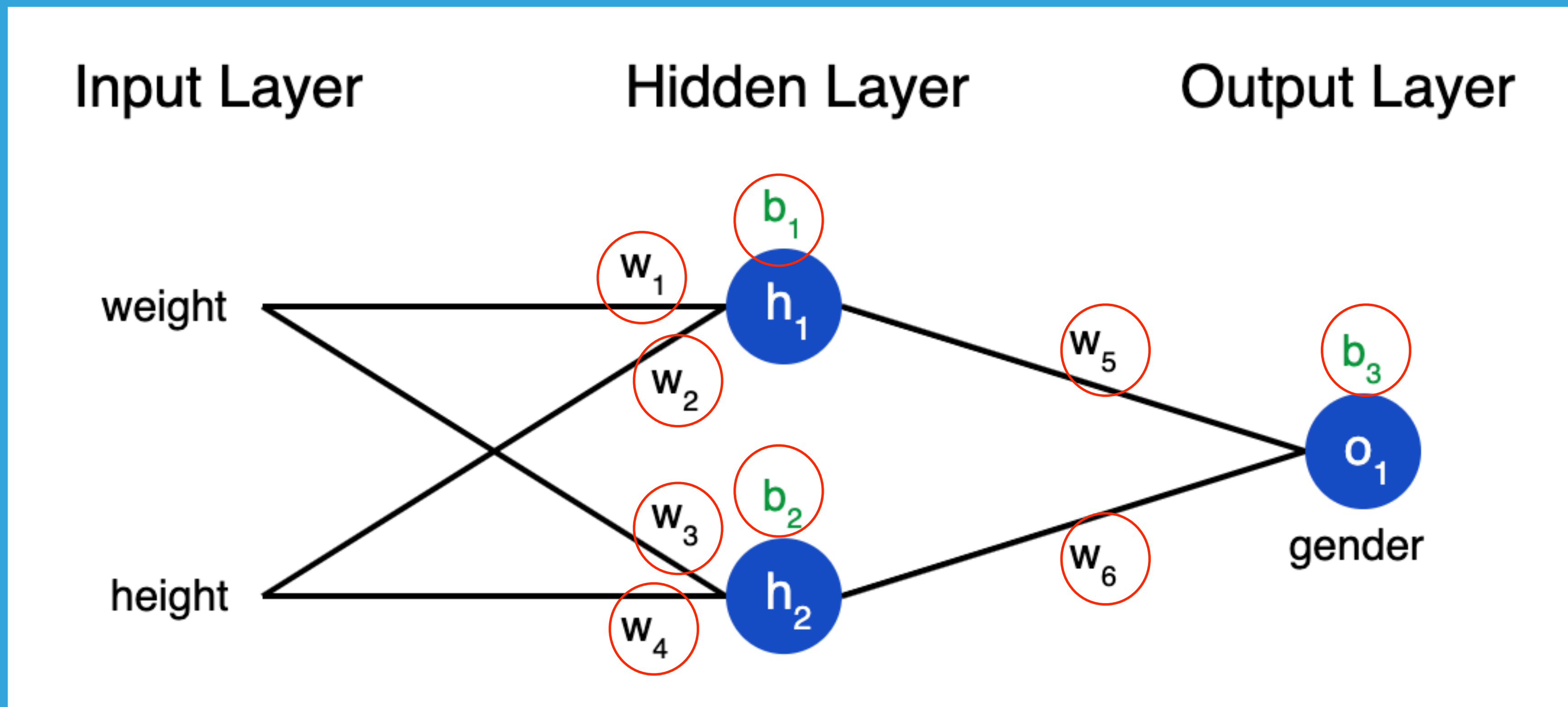$$gender = \beta_0 + \beta_1 weight + \beta_2 height$$

$$gender = \beta_0 + \beta_1 weight + \beta_2 height$$

VS

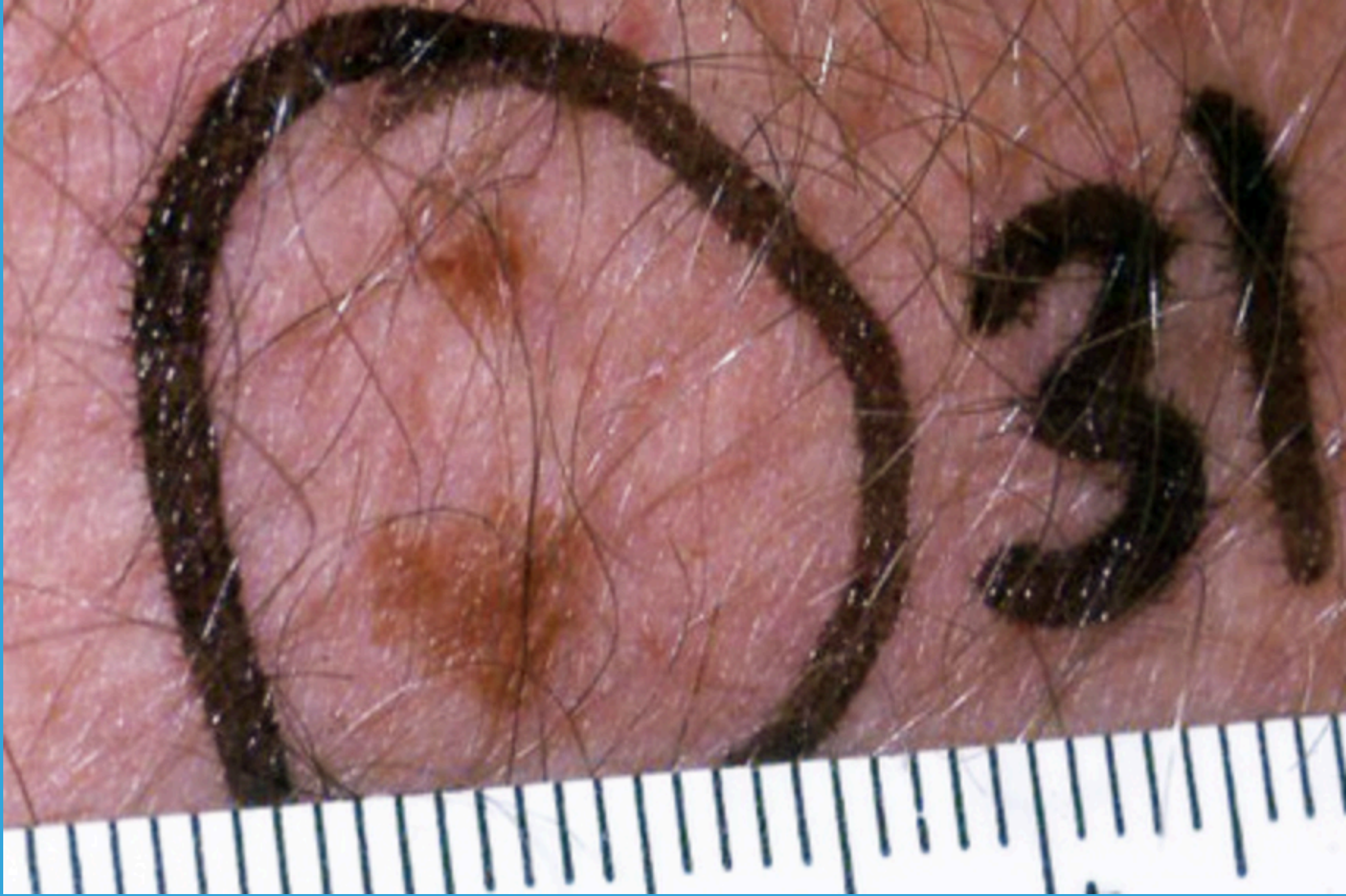$$gender = \beta_0 + \beta_1 weight + \beta_2 height$$

VS

# It's all about the Data-Generating Process

$$mpg = \beta_0 + \beta_1 cyl + disp$$

OR

$$mpg = \beta_0 + \beta_1 cyl + \beta_2 cyl^2 + \beta_3 log(disp)$$

?

# 1. INTERPRETATION MATTERS.
# 2. LINEARITY ISN'T RESTRICTIVE.
# 3. MO' SQUIGGLY = MO' OVERFITTING.
# 4. SMALL DATA'S OK.
# 5. IT'S ALL ABOUT THE DGP.

@alexgold    https://github.com/akgold/dsdc_linear_models