

Homework 2

Ryo Iwata

24 February, 2024

Homework 2

Instructions:

0. Start early and leave time to review and revise your submission well in advance of the due date.
1. Read the problem in its entirety before starting. Often, we can ignore certain variables while working on early problems and then pay attention to them in later problems. Students can waste time trying to incorporate these effects too early and with limited guidance on how to do this.
2. Plan your approach to tidying the data and preparing it for figures and analysis.
3. Translate your plan into well structured comments (in your code chunks)
4. Translate your comments into code, leaving the comments in place (so we understand your intent)
5. Ensure your code, figure and text answers are neat and tidy, grammatically correct, clear, and concise. You may find the reindent lines (`ctrl-i`) or reformat code (`ctrl-A` i.e. `ctrl-shift-a`) useful. Break arguments to functions onto different lines and ensure your code and comments do not flow off the edge of the page.
6. Review your output (pdf) to ensure your outputs render well.
7. All tables, formulae, etc., should be carefully typeset. Do not just dump outputs from functions.

Note: We will grade both your answer and how you arrived at said answer (i.e., your code). Be sure to use sensible and sufficiently descriptive tibble, variable and model names. Make efficient and appropriate use of the tidyverse functionality.

Problem #1 (60 pts):

Your laboratory is evaluating the diffusion of drug delivery vehicles (particles) injected into the brain. The lab is seeking to understand how the diameter *and* surface charge of the drug delivery vehicle affects the vehicle's diffusion.

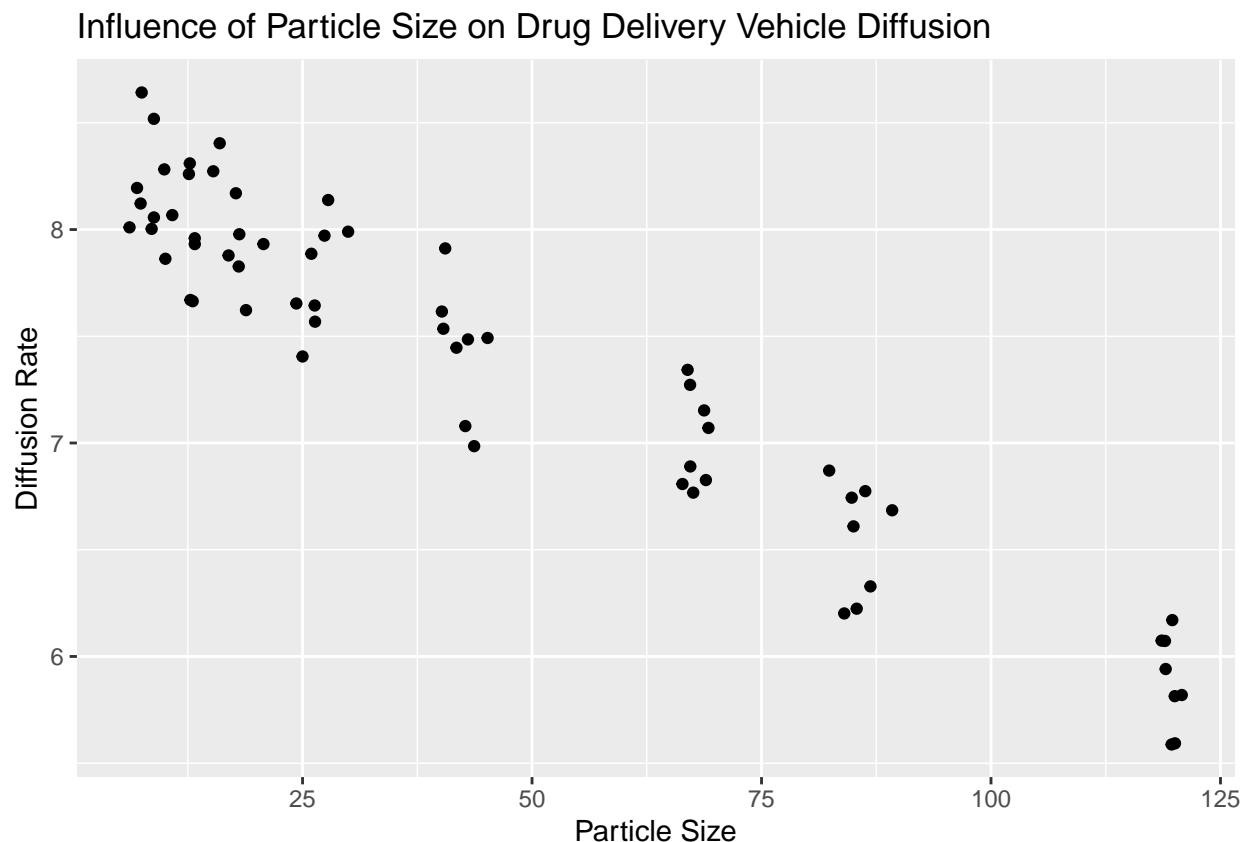
Load and evaluate the structure of data in HW2Data sheet DiffSizeCharge (only workbook/sheet in file)

```
drug_diffusion_raw <- readxl::read_excel("./data/HW2data.xlsx", sheet="DiffSizeCharge")
```

1) Investigate Particle Size and Diffusion Rate. At this point, ignore the effects of surface charge.

a) Create a plot to visualize the effect of particle diameter on particle diffusion. Think about the regression model you may be fitting, and place the appropriate variables on the x axis and y axis. Do not include any summary or smooth geometries (2 pts)

```
# Plotting scatter plot of Particle size against Diffusion rate
drug_diffusion_raw |> ggplot(aes(x=`Particle Size`, y=`Diffusion Rate`)) +
  geom_point() +
  labs(title = "Influence of Particle Size on Drug Delivery Vehicle Diffusion")
```



c) What are the model's estimated coefficients? Write out the generic formula for this model (i.e., with β terms) and then again with the estimates correctly substituted (4 pts)

```
tidy_model_size_to_diffusion <- tidy(size_to_diffusion_model)

kable(tidy_model_size_to_diffusion)
```

term	estimate	std.error	statistic	p.value
(Intercept)	8.3279	0.0489	170.23	0
Particle Size	-0.0204	0.0008	-25.18	0

```
intercept_size_to_diffusion_model <- pull(filter(tidy_model_size_to_diffusion,
                                                    term=="(Intercept)"), estimate)
slope_size_to_diffusion_model <- pull(filter(tidy_model_size_to_diffusion,
                                              term=="`Particle Size`"), estimate)

print(paste("Coefficient for intercept: ", intercept_size_to_diffusion_model))
```

```
## [1] "Coefficient for intercept: 8.32792802062674"
```

```
print(paste("Coefficient for Particle Size: ", slope_size_to_diffusion_model))
```

```
## [1] "Coefficient for Particle Size: -0.0203730765105808"
```

Answer: Generic formula for this model $\hat{y} = \beta_0 + \beta_{\text{diameter}}x_{\text{diameter}}$

Answer: Formula for this model with the estimates substituted $\hat{y} = 8.3279 + -0.0204x_{\text{diameter}}$

d) Calculate and present in a table the SS_x , SS_y , SS_{xy} , $SS_{\text{Regression}}$, and SS_{Error} for this model. (10 pts)

```
# Adding lm information to data
augment_size_to_diffusion_model <- augment(size_to_diffusion_model)
```

```
# Testing if the regression coefficients are equal to zero
tidy_anova_size_to_diffusion <- tidy(car::Anova(size_to_diffusion_model))
```

```
size_to_diffusion_SSx <- pull(summarise(augment_size_to_diffusion_model,
                                         SSx = sum((`Particle Size` - mean(`Particle Size`))^2)
                                         )
                              )
```

```
# Calculating the sum of squares of Y, the predictor variable
size_to_diffusion_SSy <- pull(summarise(augment_size_to_diffusion_model,
                                         SSy = sum((`Diffusion Rate` - mean(`Diffusion Rate`))^2)
                                         )
                              )
```

```
# Calculating the sum of squares of X and Y
size_to_diffusion_SSxy <- pull(summarise(augment_size_to_diffusion_model,
                                         SSxy = sum((`Particle Size` - mean(`Particle Size`)) * (`Diffusion Rate` - mean(
                                         )
                                         )
                              )
```

```

# Calculating the sum of squares of regression or how well the model represents the data
size_to_diffusion_SSRegression <- pull(summarise(augment_size_to_diffusion_model,
  SSRegression = sum((.fitted - mean(`Diffusion Rate`)) ^ 2)
))

# Calculating the sum of squares of error/residual
# which measures the discrepancy between the data and the model)
size_to_diffusion_SSError <- pull(summarise(augment_size_to_diffusion_model,
  SSError = sum((.fitted - `Diffusion Rate`) ^ 2)
))

kable(tibble(term = c("SSx", "SSy", "SSxy", "SSRegression", "SSError"),
  value = c(size_to_diffusion_SSx, size_to_diffusion_SSy, size_to_diffusion_SSxy, size_to_diffusion_SSRegression, size_to_diffusion_SSError)
))

```

term	value
SSx	89746.034
SSy	40.893
SSxy	-1828.403
SSRegression	37.250
SSError	3.643

```

# Calculating the sum of squares of X, the predictor variable
SSx_size_to_diffusion <- pull(summarise(augment_size_to_diffusion_model,
  SSx = sum((`Particle Size` - mean(`Particle Size`))^2)
))

# Calculating the sum of squares of error/residual
# which measures the discrepancy between the data and the model)
SSError_size_to_diffusion <- pull(filter(tidy_anova_size_to_diffusion,
  term=="Residuals"), sumsq)

# Calculating the sum of squares of regression
# or how well the model represents the data
SSRegression_size_to_diffusion <- pull(filter(tidy_anova_size_to_diffusion,
  term=="`Particle Size`"), sumsq)

# Calculating the sum of squares of Y, the predictor variable
SSy_size_to_diffusion <- SSError_size_to_diffusion + SSRegression_size_to_diffusion

# Calculating the sum of squares of X and Y
SSxy_size_to_diffusion <- pull(summarise(augment_size_to_diffusion_model,
  SSxy = sum((`Particle Size` - mean(`Particle Size`)) *
    (`Diffusion Rate` - mean(`Diffusion Rate`))
))

kable(tibble(term = c("SSx", "SSy", "SSxy", "SSRegression", "SSError"),

```

```

    value = c(SSx_size_to_diffusion, SSy_size_to_diffusion,
              SSxy_size_to_diffusion, SSregression_size_to_diffusion,
              SSerror_size_to_diffusion)
  )
)

```

term	value
SSx	89746.034
SSy	40.893
SSxy	-1828.403
SSRegression	37.250
SSerror	3.643

e) What is R^2 for this model (2 pts)?

```

# Calculating R^2, which is another measure of fit that is between 0 and 1
# Which is the proportion of the variability of Y that can be explained by X
size_to_diffusion_r_squared <- (size_to_diffusion_SSy - size_to_diffusion_SSerror) / size_to_diffusion_SSy

print(paste("R-squared for the model: ", size_to_diffusion_r_squared))

```

```
## [1] "R-squared for the model: 0.91091452224952"
```

f) Using the estimated slope term (β_1) and its associated error, use a t-test to evaluate whether the slope term is zero or non-zero. What is the t-statistic and associated p-value (2 pts)?

```

size_to_diffusion_diameter_tstat <- pull(filter(tidy_model_size_to_diffusion, term=="`Particle Size`"),
    statistic)

size_to_diffusion_diameter_pvalue <- pull(filter(tidy_model_size_to_diffusion, term=="`Particle Size`"),
    p.value)

print("t-statistic and associated p-value for slope term")

```

```
## [1] "t-statistic and associated p-value for slope term"
```

```
print(paste("T(", df.residual(size_to_diffusion_model), "):", size_to_diffusion_diameter_tstat, "p = ",
```

```
## [1] "T( 62 ): -25.1785781155024 p = 2.9349951753222e-34"
```

Answer: Assuming significance with a p-value less than 0.05, we can reject the null hypothesis that the slope term is zero. Which is evidence suggesting a significant linear relationship between the independent and dependent variables in the model

g) Does particle size have a significant effect on particle diffusion with this model? (2 pts)

```
kable(tidy_anova_size_to_diffusion)
```

term	sumsq	df	statistic	p.value
Particle Size	37.250	1	634	0
Residuals	3.643	62	NA	NA

Answer: Assuming significance with a p-value less than 0.05, we can reject the null hypothesis that all slope coefficients in the model are equal to zero because the p-value for the F-statistic is lower than the threshold. This is evidence that the model fits the data significantly better than a model without predictors.

```
kable(tidy_model_size_to_diffusion)
```

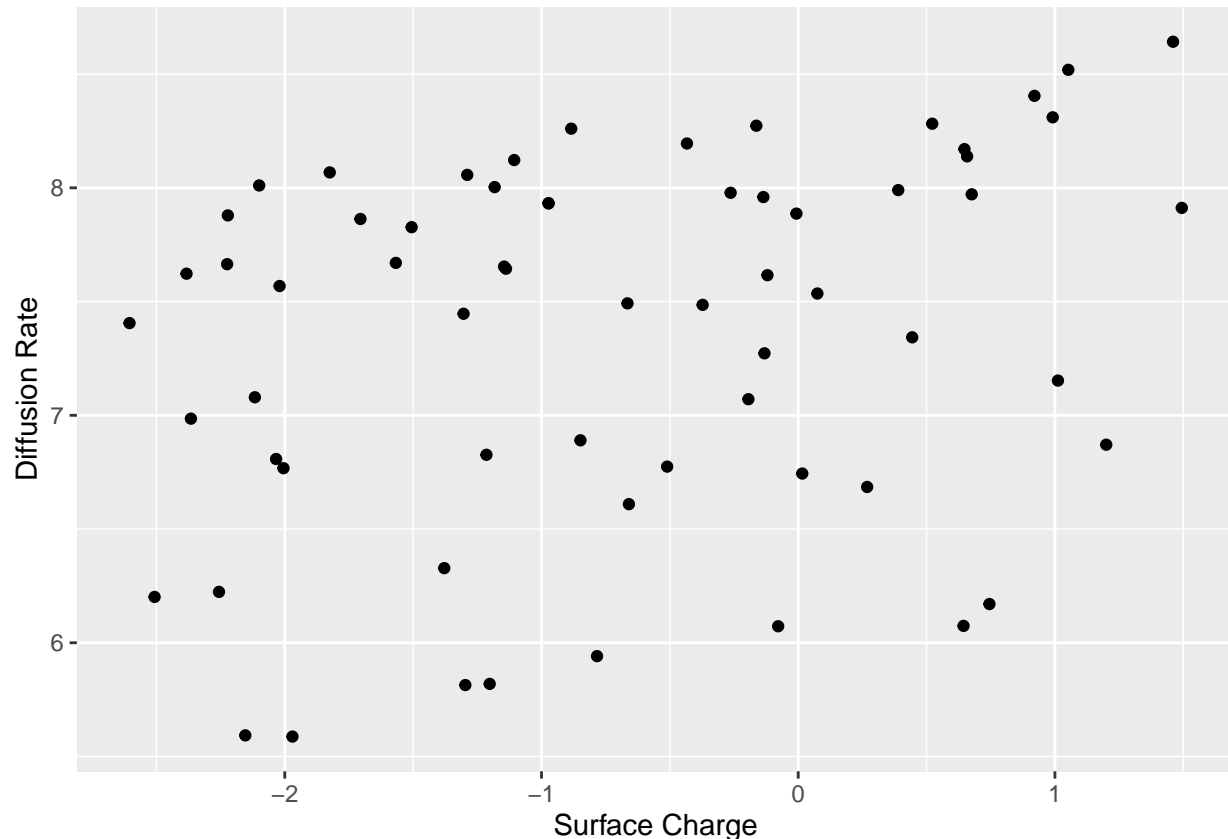
term	estimate	std.error	statistic	p.value
(Intercept)	8.3279	0.0489	170.23	0
Particle Size	-0.0204	0.0008	-25.18	0

We can also reject the null hypothesis that the coefficient associated with particle size is equal to zero because the p-value for the t-statistic is lower than the threshold. The analysis suggests that particle size has a significant effect on particle diffusion because we have evidence that there is a significant linear relationship between particle size and particle diffusion.

2) In the same experiment, these particles had different surface charges. You decide to investigate the effect of surface charge on particle diffusion.

a) Create a plot to visualize the effect of surface charge on particle diffusion in whole blood. At this point, ignore the effect of particle diameter (2 pts).

```
drug_diffusion_raw |> ggplot(aes(x=`Surface Charge`, y=`Diffusion Rate`)) +  
  geom_point()
```



b) Fit a univariate (1-variable) regression model to these data. (4 pts)

c) What are the model's estimated coefficients? Write out the formula for the model with generic terms as well as substituting the estimated coefficients (2 pts)

d) Calculate and present in a table the SS_x , SS_y , SS_{xy} , $SS_{\text{Regression}}$, and SS_{Error} for this model. (10 pts)

e) What is R^2 for this model (2 pts)?

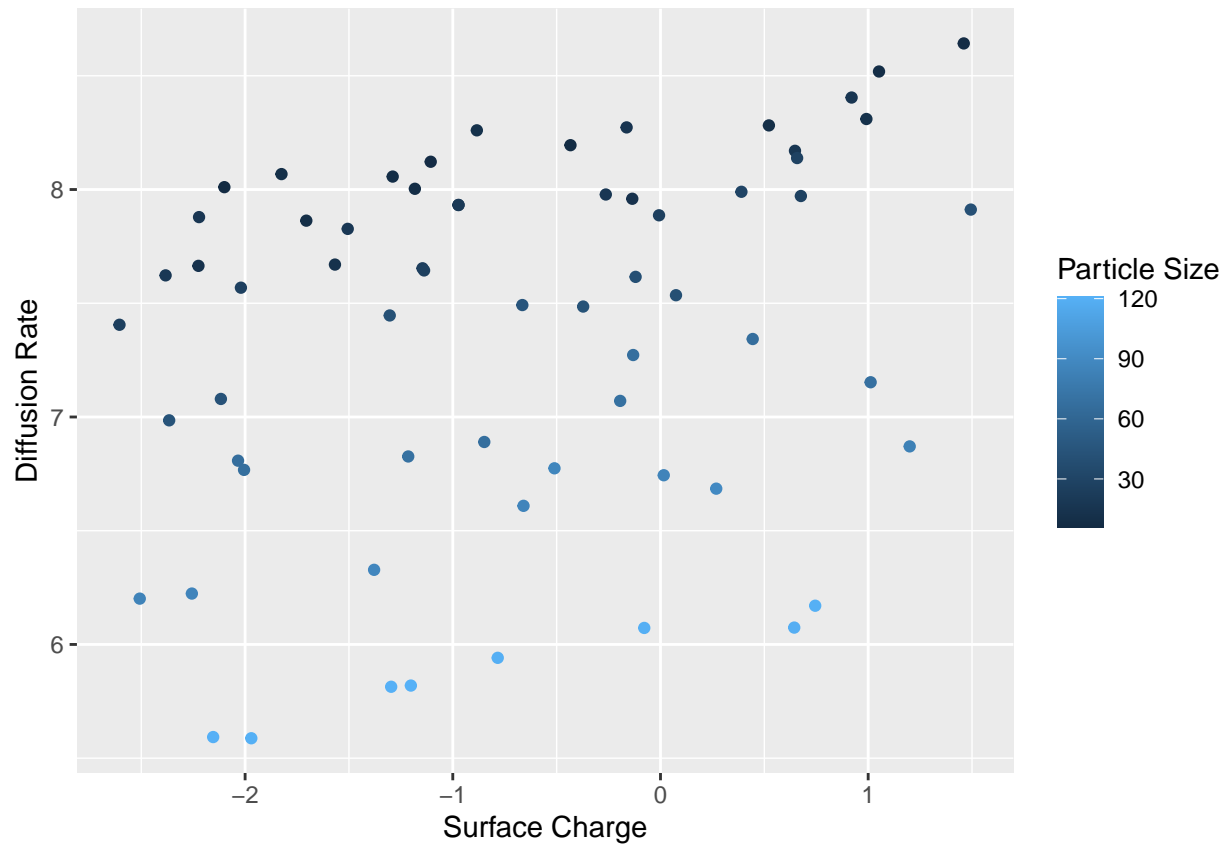
f) Using the estimated slope term (β_1) and its associated error, use a t-test to evaluate whether the slope term is zero or non-zero. What is the t-statistic and associated p-value (2 pts)?

g) Does particle size have a significant effect on particle diffusion in this model? (2 pts)

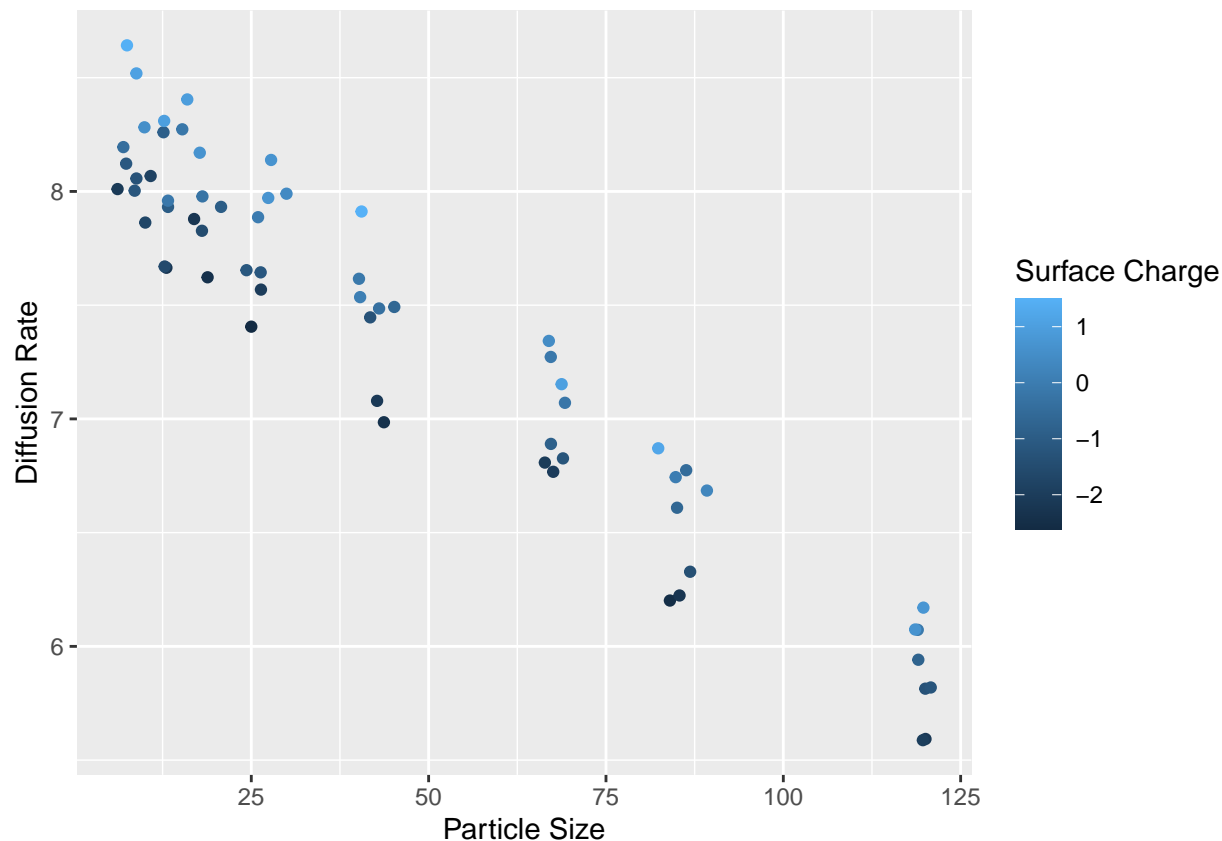
3) Multiple Regression

a) Create a 2-D scatterplot to visualize combined effect of surface charge and particle diameter on particle diffusion in whole blood. Think about the regression model you may be fitting, and use appropriate aesthetic mappings. Do not include summary or smooth geoms/stats (3 pts).

```
drug_diffusion_raw |> ggplot(aes(x=`Surface Charge`, y=`Diffusion Rate`, color=`Particle Size`)) +  
  geom_point()
```



```
drug_diffusion_raw |> ggplot(aes(x=`Particle Size`, y=`Diffusion Rate`, color=`Surface Charge`)) +  
  geom_point()
```

b) Fit a 2-variable regression model to these data. Do not include any interactions or additional terms in your model. (4 pts)

```
diameter_and_charge_model <- lm(`Diffusion Rate` ~ `Particle Size` + `Surface Charge`, drug_diffusion_r
```

c) What are the model coefficients and how do they compare to the prior models' coefficients? Write out all three model equations and evaluate. (6 pts)

```
tidy_diameter_and_charge_model <- tidy(diameter_and_charge_model)
```

```
kable(tidy_diameter_and_charge_model)
```

term	estimate	std.error	statistic	p.value
(Intercept)	8.4597	0.0212	399.21	0
Particle Size	-0.0203	0.0003	-61.83	0
Surface Charge	0.1929	0.0109	17.76	0

```
tidy_diameter_and_charge_model
```

```
## # A tibble: 3 x 5
##   term                estimate std.error statistic  p.value
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)        8.46     0.0212    399.   6.05e-106
## 2 `Particle Size`   -0.0203 0.000328  -61.8   9.75e- 57
```

```
## 3 `Surface Charge`    0.193    0.0109        17.8 8.61e- 26
diff_model_intercept <- pull(filter(tidy_diameter_and_charge_model, term=="(Intercept)"), estimate)
diff_model_part_size_slope <- pull(filter(tidy_diameter_and_charge_model, term=="`Particle Size`"), estimate)
diff_model_part_charge_slope <- pull(filter(tidy_diameter_and_charge_model, term=="`Surface Charge`"), estimate)

print(paste("Coefficient for intercept: ", diff_model_intercept))

## [1] "Coefficient for intercept: 8.45965218592573"
print(paste("Coefficient for Particle Size: ", diff_model_part_size_slope))

## [1] "Coefficient for Particle Size: -0.0203034215700633"
print(paste("Coefficient for Surface Charge: ", diff_model_part_charge_slope))

## [1] "Coefficient for Surface Charge: 0.192888719839729"
# TODO: Write out beta terms
```

d) What are $SS_{\text{ParticleSize}}$, $SS_{\text{SurfaceCharge}}$, $SS_{\text{ParticleSizeY}}$, $SS_{\text{SurfaceChargeY}}$, SS_Y , $SS_{\text{Regression}}$, and SS_{Error} (10 pts)?

e) What is R^2 for this model (2 pts)?

f) How does this model's R^2 compare to the prior two models, and why is it so? (4 pts)?

g) Use a t-test to evaluate whether each term is zero or non-zero. (4 pts)?

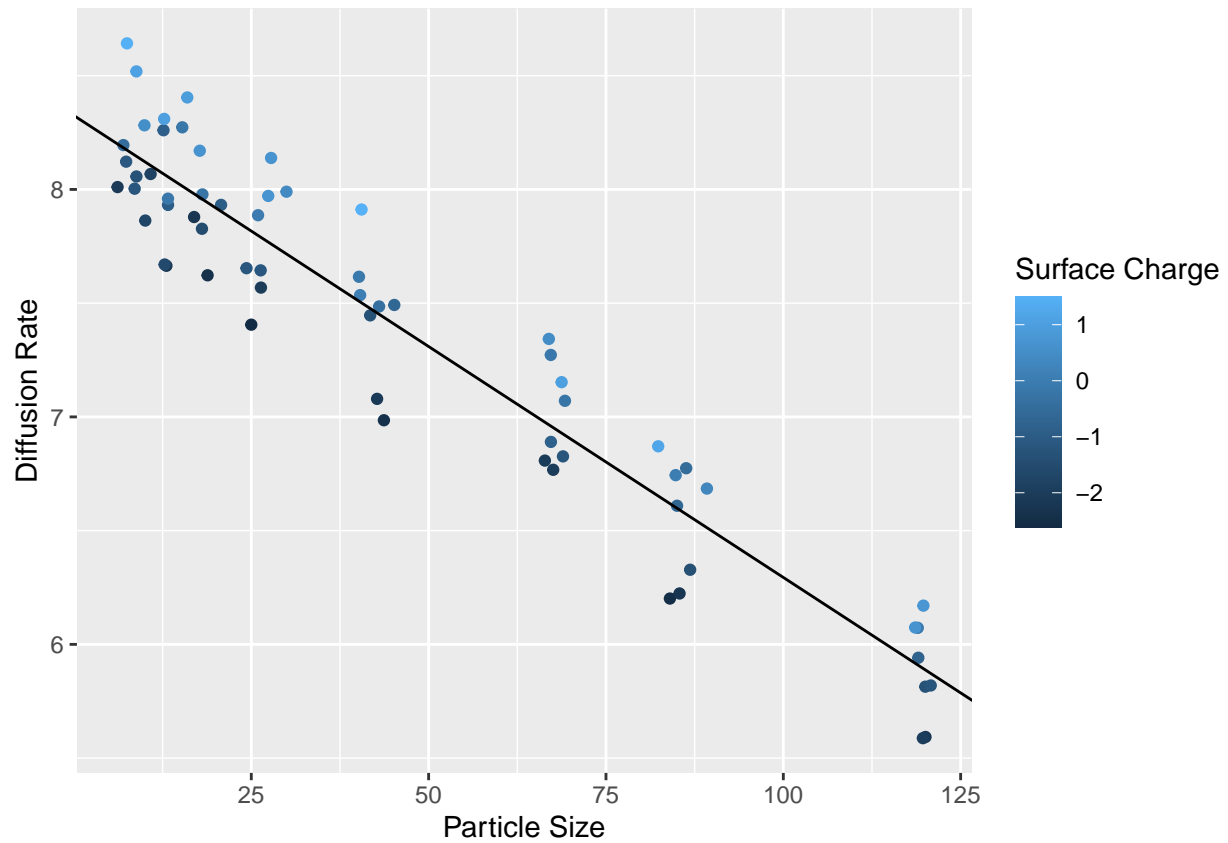
h) Does size have a significant effect on particle diffusion and, if so, what is that effect? (2 pt)

i) Does surface charge have a significant effect on particle diffusion and, if so, what is that effect? (2 pt)

j) Extend your figures (use your code for your figures from part a of this problem as a base) to show both the raw data and the model using `geom_abline`. Do not use `stat_smooth` or other related functions. Think very carefully about where your model's plane is and where/how to represent it in your figures. (8 pt)

```
mean_charge <- drug_diffusion_raw |> pull(`Surface Charge`) |>
  mean()

drug_diffusion_raw |> ggplot(aes(x=`Particle Size`, y=`Diffusion Rate`, color=`Surface Charge`)) +
  geom_point() +
  geom_abline(intercept=diff_model_intercept + diff_model_part_charge_slope * mean_charge,
             slope=diff_model_part_size_slope)
```



k) If you ran the experiment that resulted in this data, where would you put your focus for optimization of diffusion rate and why? (5 pts)

Answer: Diameter because it significantly affects diffusion rate while charge does not in a model that looks at both's effects independently.