# Cheese Products Nutrition Profiles

Jean Batista

This paper uses R to complete exploratory data analysis on the nutritional profile of cheese products. Using a dataset *cheese_data* that holds information on 60 different cheeses, categorized by their type and texture, as well as their macronutrients. I will show how to predict the calories of cheeses, using different regression models.

2025-03-22

## Reading the file

```
cheese_data <- read.csv("cheese_data.csv")

cheese_data
```

```
##                           brand                type  calories
protein
## 1                         Kraft         Cheddar(28g) 110.00000
9.000000
## 36                  Great Value        Parmesan(28g) 100.00000
```

This data has a couple of numerical variables found in the nutrition label of the products. For the values, i manually imported the data by visiting each product's website or checking on google images if the website did not have the nutrition info. Calories, protein, carbs, and fats are numerical variables. The "type" column will specify which category of cheese, such as cheddar or brie. In this dataset i've chosen 6 different types of cheeses, and 10 brands per cheese. That way i'd have a fairly large sample of data to work with, with some diversity as well. The serving size for each data sample is standardized to 1 oz or 28 grams. Note : Cottage cheese has 2 variations, one which is standardized to 28 grams, and the other which is 125 grams, the recommended serving. The reason this is important is because cottage cheese being liquid like affects the amounts of nutrients per gram, something we will see throughout this project.
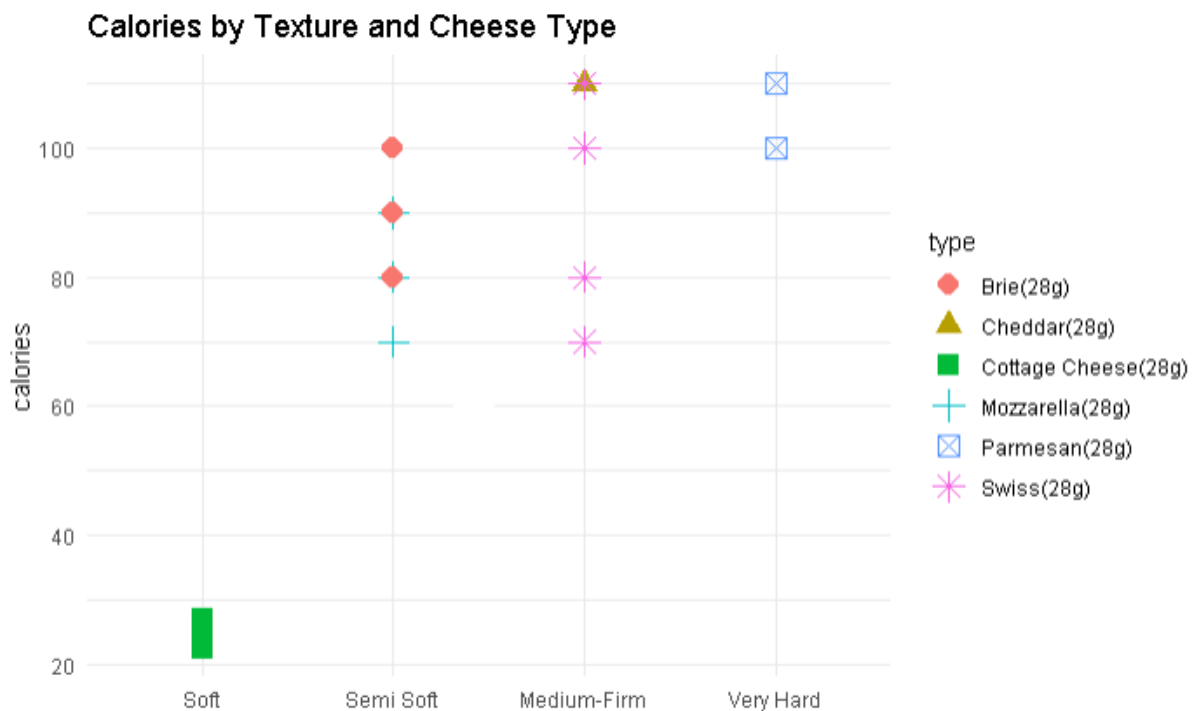
Library load

```
library(tidyverse)

library(fivethirtyeight)
library(moderndive)
library(readr)
```

Making 2 datasets, one where cottage cheese is standardized to 28 grams and one where the serving size stays at 125 grams. 28 grams will be better for identifying the characteristics that can predict calories. While it's true that it isn't reasonable to eat only 28 grams of cottage cheese, due to its water based texture, I think there is an observation to be made about why the texture affects nutrient density, we will uncover more soon.

```r
cheese_data_cot_28 <- cheese_data[c(1:50, 61:70), ]

cheese_data_cot_125 <- cheese_data[1:60,]

ggplot(cheese_data_cot_28, aes(x = Texture, y = calories, shape = type, color = type)) +
  geom_point(size = 4) +
  labs(title = "Calories by Texture and Cheese Type")
```
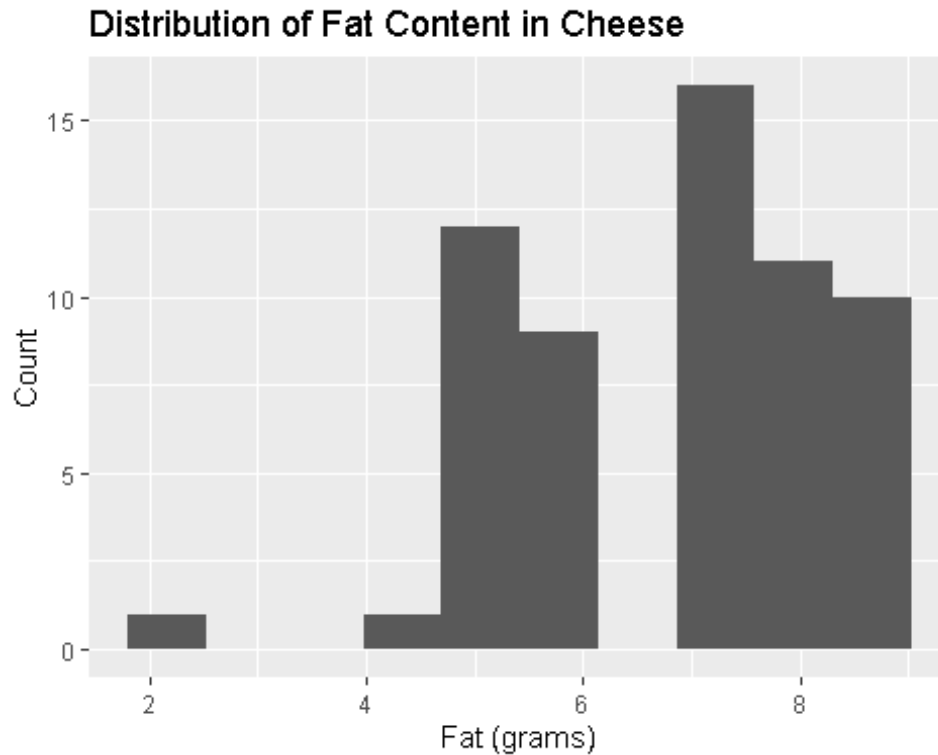
## Calories by Texture and Cheese Type



Theres a positive relationship between texture and calories. Texture gives a general expectation, but cheese type is the stronger predictor of calories. Even within the same texture category, cheeses differ significantly in caloric content, as seen by the difference in calories by cheddar and Swiss. Cottage cheese as an outlier can be studied further to evaluate why exactly it has lower nutrient density.

Distribution of Fat

```
ggplot(cheese_data_cot_125, aes(x = fat)) +
  geom_histogram(bins = 10) +
  labs(title = "Distribution of Fat Content in Cheese", x = "Fat (grams)", y
= "Count")
```

## Distribution of Fat Content in Cheese



In this graph, we were looking at all the cheeses at their intended serving size. We have some data points on the left that we want to know further about.

(https://extension.psu.edu/fat-facts-the-right-amount-for-a-healthy-diet) Research shows a gram of fat is equal to 9 calories. Where Carbs and protein are both equal to 4. Based on this information, cheese will often be a high calorie snack, considering its high fat content. Consumers looking for a healthier alternative can filter for low fat content.

```
cheese_data %>%
  filter(fat < 6)
```

```
##                        brand                type   calories  fat
## 4             Organic Valley     Mozzarella(28g)   70.00000  5
```
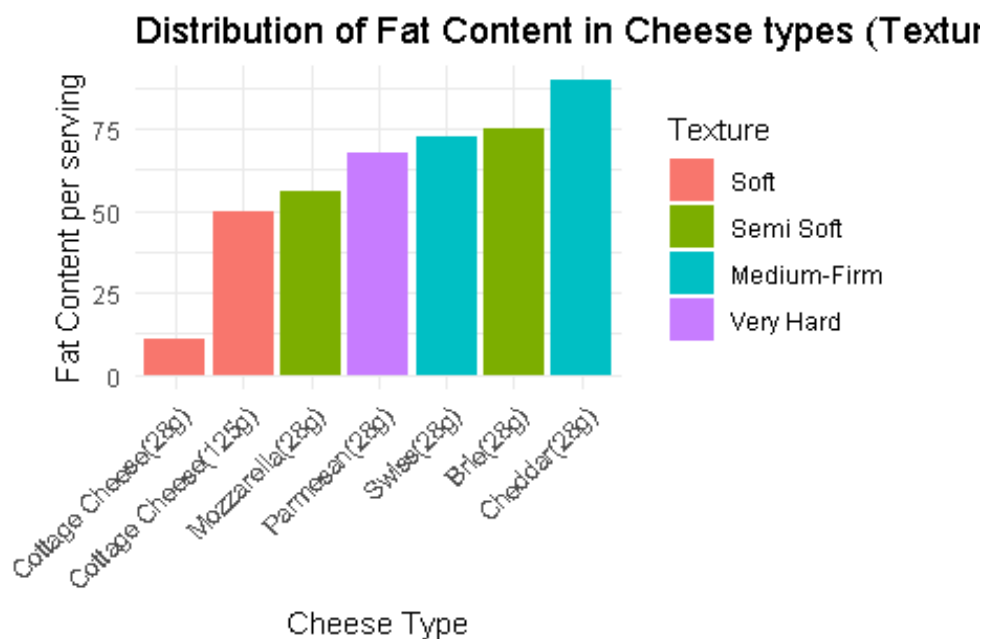
*4 more mozzarellas…*
```
## 5               Alpine Lace          Swiss(28g)   70.00000  4.5
## 6               Breakstone's  Cottage Cheese(125g) 110.00000  2.5
```

*5 more cottage cheese…*

These are cheeses with the least fat content. What's interesting is that cottage cheese has about 57% more calories than mozzarella on average per serving size, while maintaining the same fat content. Also, Alpine Lace is the only swiss cheese with low fat content. This could be that they use a different type of milk for their cheese, perhaps reduced fat milk.
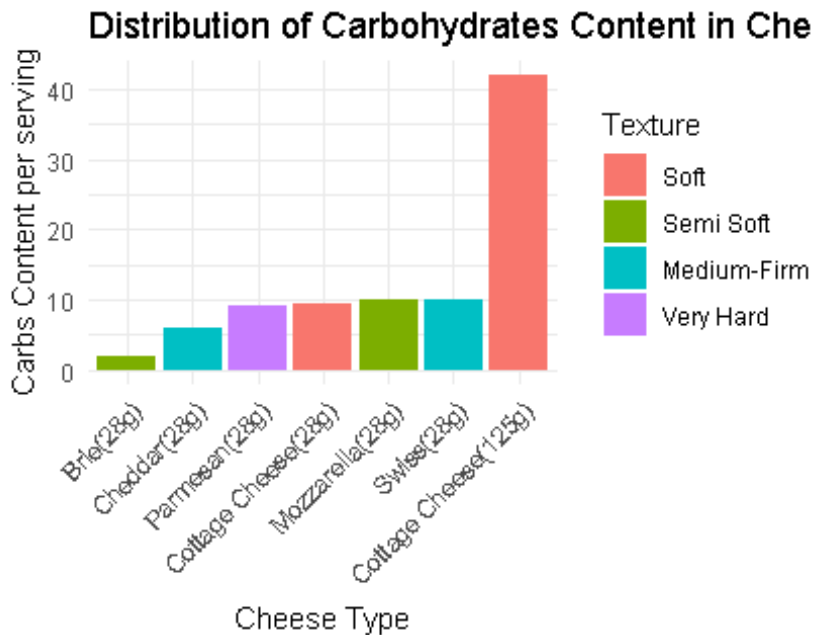
I want to see the relationship between our categorical variable Texture and some of the numerical variables, cottage cheese will be graphed with both 28 grams and 125 grams to showcase the effect of texture on nutrient density, and the need to increase serving size to accommodate for this difference.

```
ggplot(cheese_data, aes(x = reorder(type, fat), y = fat, fill = Texture)) +
  geom_col() +
  scale_fill_discrete(limits = c("Soft", "Semi Soft", "Medium-Firm", "Very
Hard")) +
  labs(title = "Distribution of Fat Content in Cheese types (Texture)",x =
"Cheese Type", y = "Fat Content per serving") +
  theme_minimal()+
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        plot.margin = margin(1, 1, 1, 1, "cm"))
```
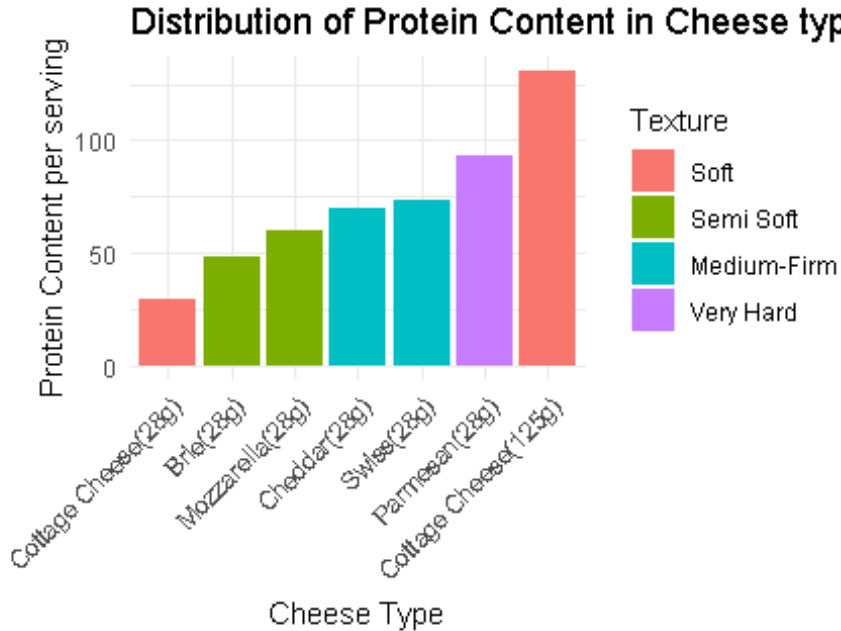


```
ggplot(cheese_data, aes(x = reorder(type, carbohydrates), y = carbohydrates,
fill = Texture)) +
  geom_col() +
  scale_fill_discrete(limits = c("Soft", "Semi Soft", "Medium-Firm", "Very
```

```
Hard")) +
  labs(title = "Distribution of Carbohydrates Content in Cheese types
(Texture)", x = "Cheese Type", y = "Carbs Content per serving") +
  theme_minimal()+
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        plot.margin = margin(1, 1, 1, 1, "cm"))
```
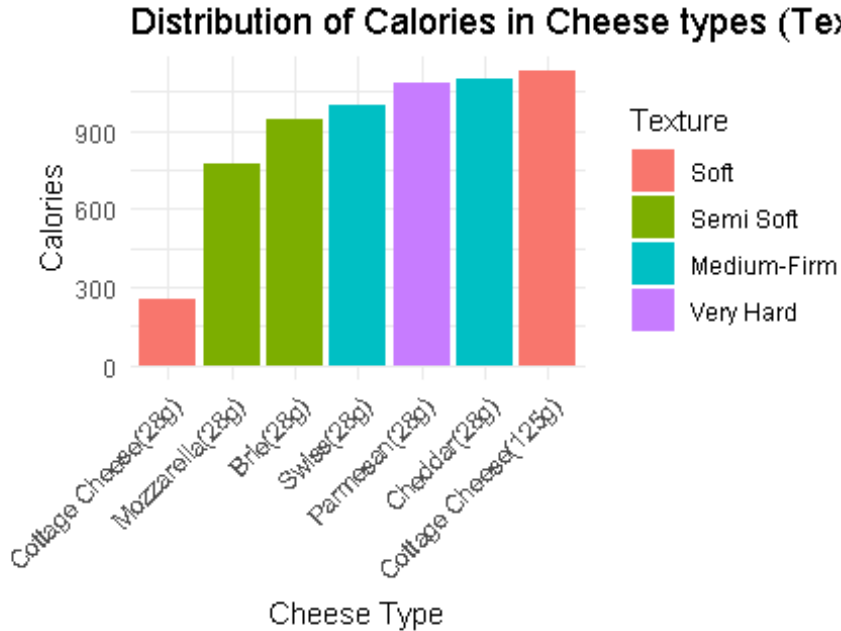


#Protein

```
ggplot(cheese_data, aes(x = reorder(type, protein), y = protein, fill =
Texture)) +
  geom_col() +
  scale_fill_discrete(limits = c("Soft", "Semi Soft", "Medium-Firm", "Very
Hard")) +
  labs(title = "Distribution of Protein Content in Cheese types (Texture)", x
= "Cheese Type", y = "Protein Content per serving") +
  theme_minimal()+
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        plot.margin = margin(1, 1, 1, 1, "cm"))
```

**Distribution of Protein Content in Cheese types (T**

Cottage cheese is the softest (it has very high moisture). Keep in mind the water content in this type of cheese makes it less nutrient dense per ounce. People would not normally eat 28 grams of cottage cheese since it would only be about a spoonful, versus the denser cheeses that are compact with nutrients.

When looking at per 28 grams, cottage cheese contains significantly less fat and protein, while competing with the other cheese types in carbohydrates.

```
ggplot(cheese_data, aes(x = reorder(type, calories), y = calories, fill =
Texture)) +
  geom_col() +
  scale_fill_discrete(limits = c("Soft", "Semi Soft", "Medium-Firm", "Very
Hard")) +
  labs(title = "Distribution of Calories in Cheese types (Texture)",x =
"Cheese Type", y = "Calories") +
    theme_minimal()+
    theme(axis.text.x = element_text(angle = 45, hjust = 1),
    plot.margin = margin(1, 1, 1, 1, "cm"))
```

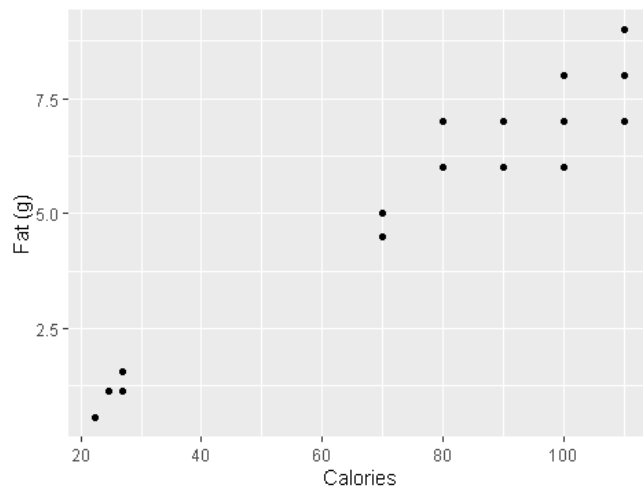## Distribution of Calories in Cheese types (Texture)



Because of this you will not expect a lot of calories per 28g of this cheese, since its fat content is so low. It is less satiating; therefore, the recommended serving is increased to 125 grams. Going forward we will only look at 28-gram servings because this insight on texture is valuable and can be evaluated further.

The trend shows as a cheese gets harder it will be more nutrient dense, containing more fat and calories. But this is a trend not absolute truth, as parmesean is the hardest cheese but contains less fat content and calories than several cheeses.

The following scatterplots will be used to visualize the correlation between calories and each numerical variable:

```
ggplot(cheese_data_cot_28, aes(x = calories, y = fat)) +
  geom_point() +
  labs(
    title = "Scatter Plot of Fat vs. Calories",
    x = "Calories",
    y = "Fat (g)"
  )
```
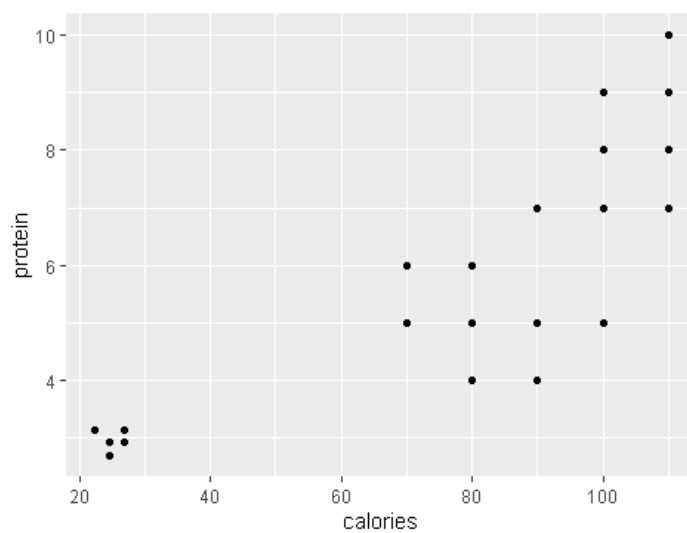
**Scatter Plot of Fat vs. Calories**



```
ggplot(cheese_data_cot_28, aes(x = calories, y = protein)) +
  geom_point() +
  labs(
    title = "Scatter Plot of Protein vs. Calories",
  )
```
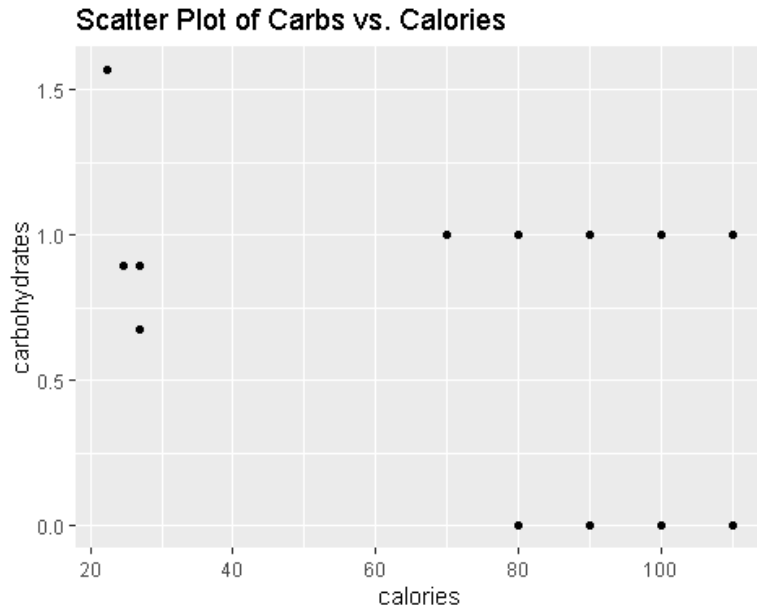
There is an upwards trend for fat.

**Scatter Plot of Protein vs. Calories**



```
ggplot(cheese_data_cot_28, aes(x = calories, y = carbohydrates)) +
  geom_point() +
  labs(
    title = "Scatter Plot of Carbs vs. Calories",
  )
```

There is an upwards trend for protein.

## Scatter Plot of Carbs vs. Calories



 There is an upwards trend for both fat and protein, but there is not any specific trend for carbohydrates. As cheese becomes richer, it will increase in fat and protein, causing the calories to increase.

We can isolate different variables to see their effect on calories. Based on the observations above we see some correlation between fat, protein, and texture. But our models below will be used to predict the expected amount.

```
# Model 1: Predict Calories from Fat
model_fat <- lm(calories ~ fat, data = cheese_data_cot_28)

# Model 2: Predict Calories from Protein
model_protein <- lm(calories ~ protein, data = cheese_data_cot_28)

# Model 3: Predict Calories from Carbs
model_carbs <- lm(calories ~ carbohydrates, data = cheese_data_cot_28)

summary(model_fat)

##
## Call:
## lm(formula = calories ~ fat, data = cheese_data_cot_28)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -14.622  -5.886  -3.357   4.114  16.642
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)   15.7730      2.7691    5.696  4.3e-07 ***
## fat           11.2641      0.4123   27.318  < 2e-16 ***
## ---
## Signif. codes:   0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.166 on 58 degrees of freedom
## Multiple R-squared:  0.9279, Adjusted R-squared:  0.9266
## F-statistic: 746.3 on 1 and 58 DF,  p-value: < 2.2e-16
```

For our Fat Model:

B0 = 15.773 : When fat is 0g we can expect 15.773 calories

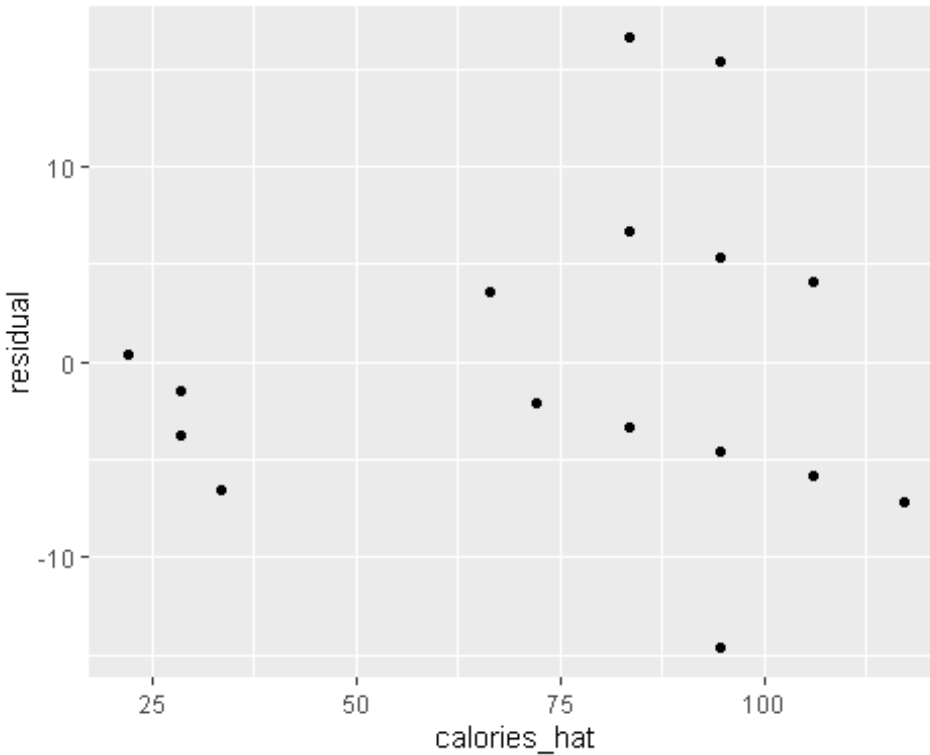B1 = 11.2641: For every gram of fat, there is an expected increase of 11.26 calories

R-squared = .9279: About 92.8% of the variation in calories is explained by Fat

RSE = 8.166: Expected vs Actual calories may differ at about 8 calories.

Fat is a strong predictor, the high R-squared suggests as such. I found it interesting that the slope is 11.26, because as we mentioned earlier, scientific research states each gram of fat is equal to 9 calories. What this tells us is that as fat increases, theres likely an increase in carbs and/or protein.

Does it match LINE?

```
get_regression_points(model_fat) -> fat_residual_info
ggplot(fat_residual_info, aes(x = calories_hat, y=residual )) +
  geom_point()
```

The plot looks randomly centered around zero, it matches with LINE.

```
summary(model_protein)

##
## Call:
## lm(formula = calories ~ protein, data = cheese_data_cot_28)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -26.368 -14.479  -3.054  14.970  30.899
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.196      6.958   1.609    0.113
## protein       11.976      1.061  11.287 2.89e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 17.01 on 58 degrees of freedom
## Multiple R-squared:  0.6872, Adjusted R-squared:  0.6818
## F-statistic: 127.4 on 1 and 58 DF,  p-value: 2.888e-16
```

For our Protein Model:

B0 = 11.196 : When protein is 0g we can expect 11.196 calories (very close to our fat slope)

B1 = 11.976: For every gram of protein, there is an expected increase of 11.976 calories
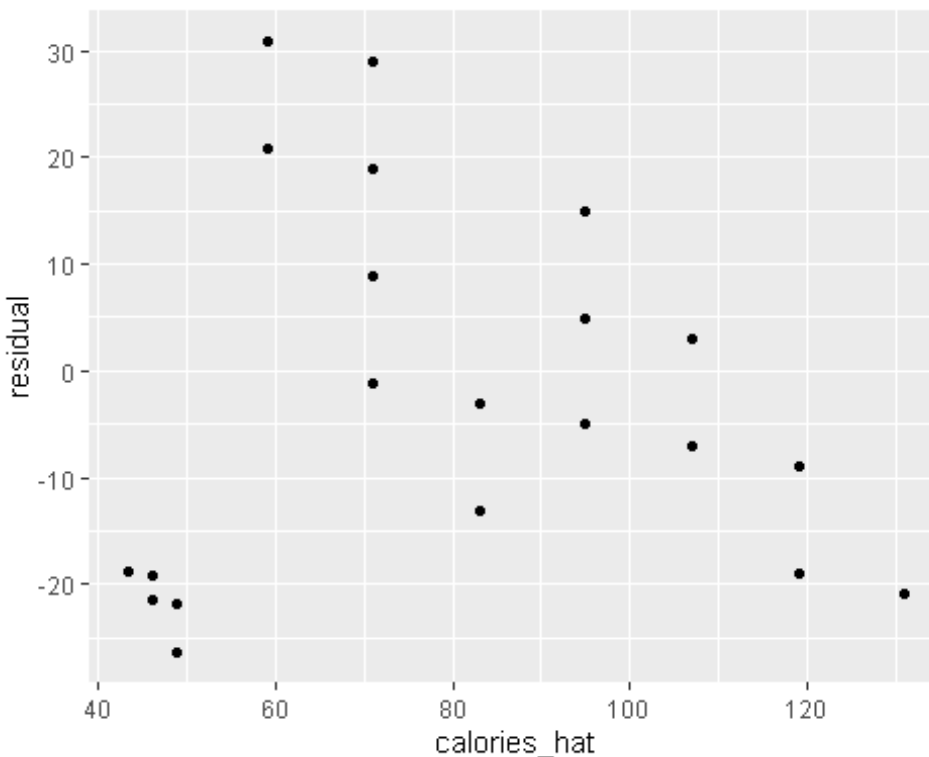
R-squared = .68772: About 68.7% of the variation in calories is explained by protein

RSE = 17.01: Expected vs Actual calories may differ at about 17 calories.

Protein is a fair predictor. Scientific research states each gram of protein is equal to 4 calories. The other calories occur because an increase of protein leads to an increase of fat, which as seen before was the primary diver of calorie variation.

Does it match LINE?

```
get_regression_points(model_protein) -> protein_residual_info
ggplot(protein_residual_info, aes(x = calories_hat, y=residual )) +
  geom_point()
```



The points show a downwards trend. The model does not match LINE.

```
summary(model_carbs)
```

```
##
## Call:
## lm(formula = calories ~ carbohydrates, data = cheese_data_cot_28)
##
## Residuals:
##     Min      1Q Median      3Q     Max
## -60.20 -12.63   10.54   27.37   27.37
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)     96.296      8.120  11.860   <2e-16 ***
## carbohydrates  -13.668      9.238  -1.479    0.144
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 29.85 on 58 degrees of freedom
## Multiple R-squared:  0.03637,    Adjusted R-squared:  0.01975
## F-statistic: 2.189 on 1 and 58 DF,  p-value: 0.1444
```

For our Carbs Model:

B0 = 96.296 : When carbs is 0g we can expect 96 calories.

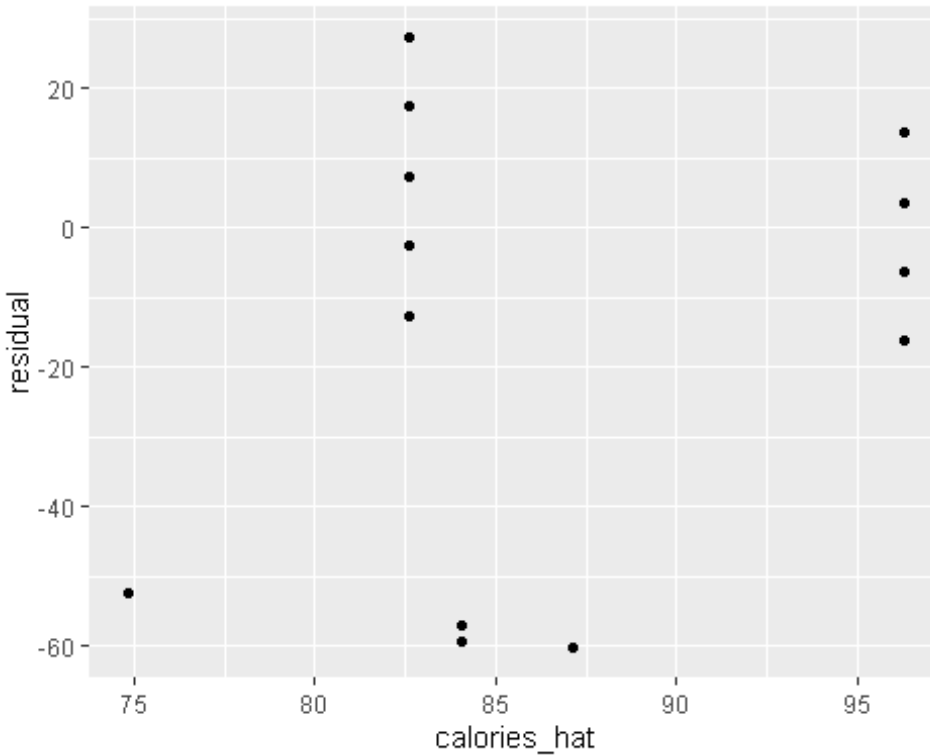B1 = -13.668: For every gram of carbs, there is an expected decrease of 13.6 calories

R-squared = .03637: About 3.6% of the variation in calories is explained by carbs

RSE = 29.85: Expected vs Actual calories may differ at about 29 calories.

Carbs have a weak relationship with cheese in this model, as suggested by the very low r squared. Scientific research states each gram of carbohydrates is equal to 4 calories. But here there is a negative slope for carbs. Due to this variable being insignificant we can't conclude there's a real negative association in this model.

Does it match LINE?

```
get_regression_points(model_carbs) -> carbs_residual_info
ggplot(carbs_residual_info, aes(x = calories_hat, y=residual )) +
  geom_point()
```

There is a vertical line pattern, and the residuals are not centered around zero, so this does not match LINE.

Since fat is our primery driver for calories, im going to use it as the variable to make predictions...

**Calories = 15.77 +11.26 (fat)**

Suppose we have cheese with 9 grams of fat, we plug in at 9 * 11.26 and get our result of 117.11

Predicted calories of 117.11, if we use cheddar with 9 grams of fat, actual calories are around 110. Meaning this model is somewhat accurate but is not exactly correct yet.

Now im going to look at the model for Fat + Protein.

```
model_fat_and_protein <- lm(calories ~ fat + protein, data =
cheese_data_cot_28)
summary(model_fat_and_protein)

##
## Call:
## lm(formula = calories ~ fat + protein, data = cheese_data_cot_28)
##
```

```
## Residuals:
##    Min     1Q Median     3Q    Max
## -4.952 -2.834  0.012  2.889  8.124
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   2.0580     1.3465   1.528    0.132
## fat           8.6558     0.2206  39.240   <2e-16 ***
## protein       4.8069     0.2725  17.638   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.242 on 57 degrees of freedom
## Multiple R-squared:  0.9888, Adjusted R-squared:  0.9884
## F-statistic:  2524 on 2 and 57 DF,  p-value: < 2.2e-16
```

For our Fat and Protein Model:

B0 = 2.05 : When carbs and protein is 0g we can expect 2.05 calories (carbs).

B1 = 8.6558: For every gram of fat, there is an expected decrease of 8.6 calories

B2 = 4.8069: For every gram of protein, there is an expected decrease of 4.8 calories

R-squared = .9888: About 98.9% of the variation in calories is explained by carbs

RSE = 3.242: Expected vs Actual calories may differ at about 29 calories.

This model is our strongest one yet, as suggested by the very high r squared. The slopes in this model are more in line with the scientific research of 9 cals per gram of fat and 4 cals per gram of protein. But here there is a negative slope for carbs.

The equation for this model is:

**Calories = 2.06 + 8.66(fat) + 4.81(protein)**

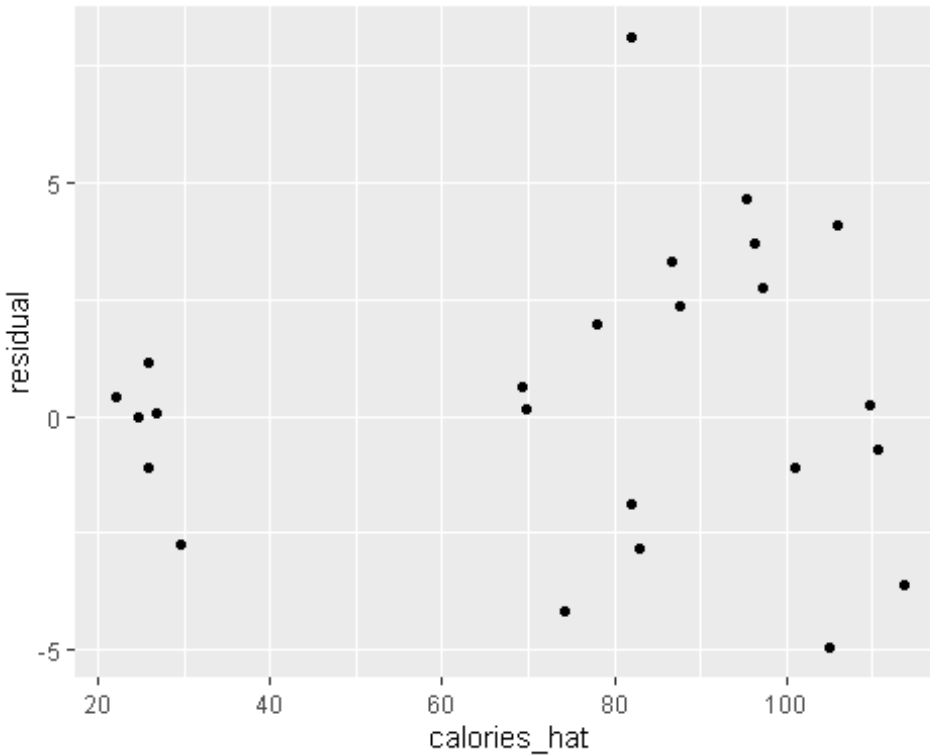Therefore 8 grams of fat and 5 grams of protein (were taking Brie cheese as an example):

( 8* 8.66 = 69.28 ) ( 5* 4.81 = 24.05 ), The sum of these 2 is 93.33. Plus the interception (2.06) equals 95.39.

Observed values of Brie are at around 100 calories. The prediction is accurate although it still could improve.

Does it match LINE?

#Look at the residual plot for this model. How well does it match LINE?

```
get_regression_points(model_fat_and_protein) -> fat_protein_residual_info
ggplot(fat_protein_residual_info, aes(x = calories_hat, y=residual )) +
  geom_point()
```

This model matches with LINE as there is no pattern and spread is centered around zero.

Now I'm going to evaluate the parallel model and interaction model, and determine which model fits best.

```
model_parallel_fat_and_type <- lm(calories ~ fat + type, data =
cheese_data_cot_28)
summary(model_parallel_fat_and_type)

##
## Call:
## lm(formula = calories ~ fat + type, data = cheese_data_cot_28)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -8.4527 -0.2541  0.0000  1.4432  8.4919
##
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)             9.4728     4.3354   2.185  0.03333 *
## fat                    11.2703     0.5682  19.834  < 2e-16 ***
## typeCheddar(28g)       -0.9054     1.4123  -0.641  0.52422
## typeCottage Cheese(28g)  3.3550     3.8016   0.883  0.38148
```

```
## typeMozzarella(28g)       4.4136     1.5601   2.829  0.00658 **
## typeParmesan(28g)        21.8892     1.1943  18.328   < 2e-16 ***
## typeSwiss(28g)            8.8176     1.1351   7.768 2.63e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.518 on 53 degrees of freedom
## Multiple R-squared:  0.9937, Adjusted R-squared:  0.993
## F-statistic:  1401 on 6 and 53 DF,  p-value: < 2.2e-16
```

This model fits the data exceptionally well, as seen in the r-squared of .9937. Not all cheese types are significant ( cheddar and cottage cheese ). But the other 4 have low p values and are significant.

The equation is:

**Calories = 11.27(FAT) + 9.47(intercept) -0.90(cheddar) + 3.35(cottage cheese) + 4.41(mozzarella) + 21.88(parmesean) + 8.81(swiss)**

I want to look at cottage cheese alone, because of the observations made earlier:

**Cottage Cheese Calories = 11.27(1.1) + 9.47(b0) + 3.35(1) = 25.217**

Our observed value for cottage cheese is at around 25.62, so this prediction came extremely close.

If I look at the prediction for cheddar calories:

**Cheddar Calories = 11.27(9) + 9.47(b0) - 8.1(1) = 102.8**

Our observed value for cheddar is at around 110, so this prediction is off by around 7.2 calories.

Now for the interaction model:

```
model_parallel_fat_and_type <- lm(calories ~ fat * type, data =
cheese_data_cot_28)
summary(model_parallel_fat_and_type)

##
## Call:
## lm(formula = calories ~ fat * type, data = cheese_data_cot_28)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -8.6139 -0.1042  0.0000  1.3861  8.3333
```

```
## 
## Coefficients: (1 not defined because of singularities)
##                            Estimate Std. Error t value Pr(>|t|)
## (Intercept)                  4.0000    11.9181   0.336   0.7386
## fat                         12.0000     1.5856   7.568 8.78e-10 ***
## typeCheddar(28g)            -2.0000     2.6294  -0.761   0.4505
## typeCottage Cheese(28g)     16.2890    12.5591   1.297   0.2007
## typeMozzarella(28g)          7.6667    14.9932   0.511   0.6114
## typeParmesan(28g)           36.0000    18.0086   1.999   0.0512 .
## typeSwiss(28g)              12.7327    12.9938   0.980   0.3319
## fat:typeCheddar(28g)            NA         NA      NA       NA
## fat:typeCottage Cheese(28g) -7.4523     3.8394  -1.941   0.0580 .
## fat:typeMozzarella(28g)     -0.3333     2.2656  -0.147   0.8836
## fat:typeParmesan(28g)       -2.0000     2.5381  -0.788   0.4345
## fat:typeSwiss(28g)          -0.5149     1.7355  -0.297   0.7680
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 2.507 on 49 degrees of freedom
## Multiple R-squared:  0.9943, Adjusted R-squared:  0.9931
## F-statistic: 848.5 on 10 and 49 DF,  p-value: < 2.2e-16
```

This model fits the data exceptionally well, as seen in the r-squared of .9943, slightly higer than the parallel model. But most of the interaction terms are insignificant, except for the cottage cheese (.0580). This coincides with our observations earlier, that cottage cheese is very soft and low in fat( or has more water content), making its calorie density lower than that of the other cheeses.

Let's look in this model in action:

**cottage cheese calories = (4 + 16.28) + (12-7.4523) fat =**

**cottage cheese calories =  20.29 + 4.55(fat)**

If we plug in 1.1 for fat,  the calories are 25.295

Our observed value for cottage cheese is at around 25.62, so this prediction came extremely close.

If we look at a type of cheese that is insignificant in this model for example Parmesan:

**calories = (4 + 36) + (12-2)fat =**

**calories = 40 + 10(fat) =**

If fat is 7, calories are 110. Our observed values are 110, so it perfectly matched.

Let's see for Mozzarella:

$$calories = (4 + 7.667) + (12 - .0333)fat =$$

$$calories = 11.667 + 11.667(fat) =$$

if fat is 6, calories are 81.668. Our observed value is 80, making this model very accurate.

Based on my testing, the interaction model is the best fit. Overall, while the interaction model is slightly more complex than the parallel model, it highlights the unique behavior of Cottage Cheese. For the other types, the effect of fat on calories is similar, but for Cottage Cheese, the lower slope in the interaction model(-7.45) captures its distinctive nutritional profile.