

Nutrition

Hesham Al Kayed

8/13/2019

1 Introduction

This project aims to build a model that can predict which food is **Raw** and which is **Processed** by examining the food nutrition values and use them as predictors. We will use Nutrition data from U.S. DEPARTMENT OF AGRICULTURE website.¹.

In **Section 2.1** we Build our data set and format the result as needed. In **Section 2.2** we explore our data to find the difference between **Raw** and **Processed** food, this exploration will be based on **food groups**, which will help us make sense of the data at hand. In **Section 3** we will build a tree based model and compare the performance of different models.

2 Exploratory data analysis

2.1 Data Structure

Downloaded data will contain a documentation file and data files. as mentioned in the documentation, data files are carrot '^' delimited. We will use and connect NUT_DATA, NUTR_DEF, FOOD_DESC and FD_GROUP tables as described in the documentation. then we will select needed fields and rename them as follow **ProdID = NDB_No** , **NutrVal = Nutr_Val**, **ShrtDesc = Shrt_Desc** , **FoodGroup = FdGrp_Desc**. (left of the equal sign is the new name)

Table 1: Data Set

ProdID	NutrVal	Units	NutrDesc	ShrtDesc	ManufacName	FoodGroup
01001	717	kcal	Energy	BUTTER,WITH SALT	NA	Dairy and Egg Products
01001	0	mg	Caffeine	BUTTER,WITH SALT	NA	Dairy and Egg Products
01001	0	mg	Theobromine	BUTTER,WITH SALT	NA	Dairy and Egg Products
01001	24	mg	Calcium, Ca	BUTTER,WITH SALT	NA	Dairy and Egg Products
01001	2	mg	Magnesium, Mg	BUTTER,WITH SALT	NA	Dairy and Egg Products
01001	24	mg	Phosphorus, P	BUTTER,WITH SALT	NA	Dairy and Egg Products

Table 2: Types and Missing Data

	ProdID	NutrVal	Units	NutrDesc	ShrtDesc	ManufacName	FoodGroup
Type	character	double	character	character	character	character	character
NA	0	0	0	0	0	426413	0
Empty	0	0	0	0	0	0	0

Table 3: Unique Values

ProdID	NutrVal	Units	NutrDesc	ShrtDesc	ManufacName	FoodGroup
7793	17446	5	97	7790	139	25

As we can see in *Table.2* all variables are complete except for **ManufacName**, which as per downloaded documentation is only available if possible. Our data set have 7 variables and 468669 records. we can see from *Table.3* that we have 7793 products and 7790 descriptions which indicates a duplicated product ID.

¹https://www.ars.usda.gov/ARSUserFiles/80400525/Data/SR-Legacy/SR-Leg_ASC.zip

Table 4: Duplicated IDs

ProdID	ShrtDesc
04657	OIL,INDUSTRIAL,PALM KERNEL (HYDROGENATED),CONFECTION FAT
04658	OIL,INDUSTRIAL,PALM KERNEL (HYDROGENATED),CONFECTION FAT
13351	BEEF,CHUCK,UNDER BLADE CNTR STEAK,BNLESS,DENVER CUT,LN,0" FA
13352	BEEF,CHUCK,UNDER BLADE CNTR STEAK,BNLESS,DENVER CUT,LN,0" FA
25000	POPCORN,OIL-POPPED,LOFAT
25001	POPCORN,OIL-POPPED,LOFAT

by removing duplicated IDs we get the below table.

Table 5: Unique Values

ProdID	NutrVal	Units	NutrDesc	ShrtDesc	ManufacName	FoodGroup
7790	17443	5	97	7790	139	25

2.2 Data Description

As we want to build a model to classify *raw foods* from *processed food* we will start by exploring the difference between raw and processed food. we will separate raw foods from processed foods by searching for **‘RAW’**, in the short description field using the following Reg. expression `'(?![[:alpha:]]RAW(?![[:alpha:]])'`, and we will consider every product in **Sausages and Luncheon Meats** as processed product and every product in **Spices and Herbs** as raw food, Available data have 6374 Processed products and 1416 raw foods.²

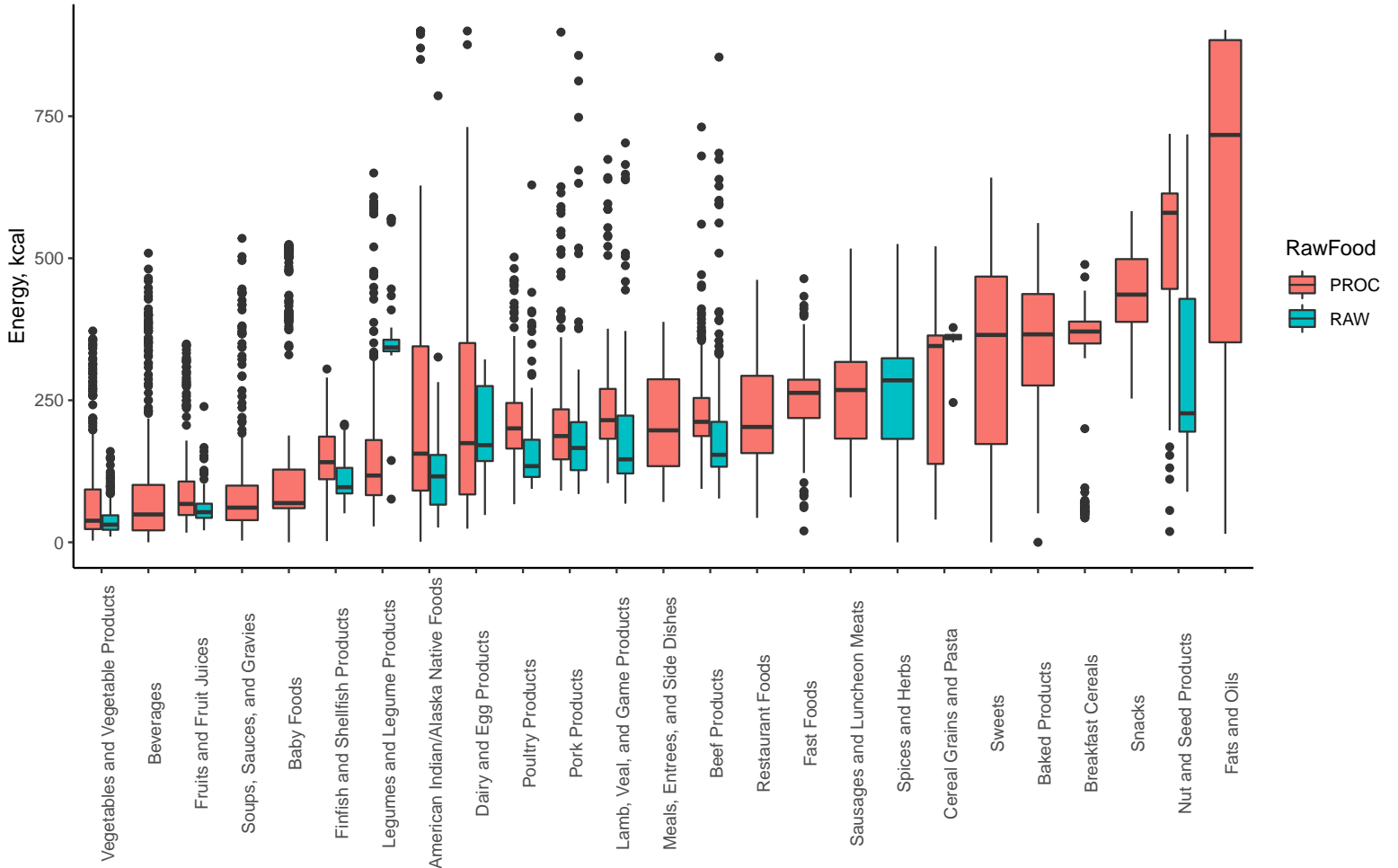


Figure 1: Comparing Raw and Processed Foods Energy Values

²We have to note that this method of differentiating between raw and processed food may not include all cases. but because we don't know how the data is constructed and there is no mention to this information in the documentation, we will assume that this method is correct

Fig.1 shows that Processed food on average have more kaloris per 100g. we can see that this does not apply to Legumes and Legume Product. Fig.2 Shows the median value for each of the energy sources in food, we can see that most raw foods do not have sugar in them and processed food does, raw Legumes and Legume Product have more protein and more carbohydrate than processed Legumes and Legume Product, which can explain why raw Legumes and Legume Product has more energy than processed Legumes and Legume Product and other food groups don't.

Fiber content per 100g does not deffer much between processed and raw foods, except for Legumes and Legume Products as shown in Fig.3. we can see in Fig.4 that all raw foods except Legumes and Legume Product have more water content.

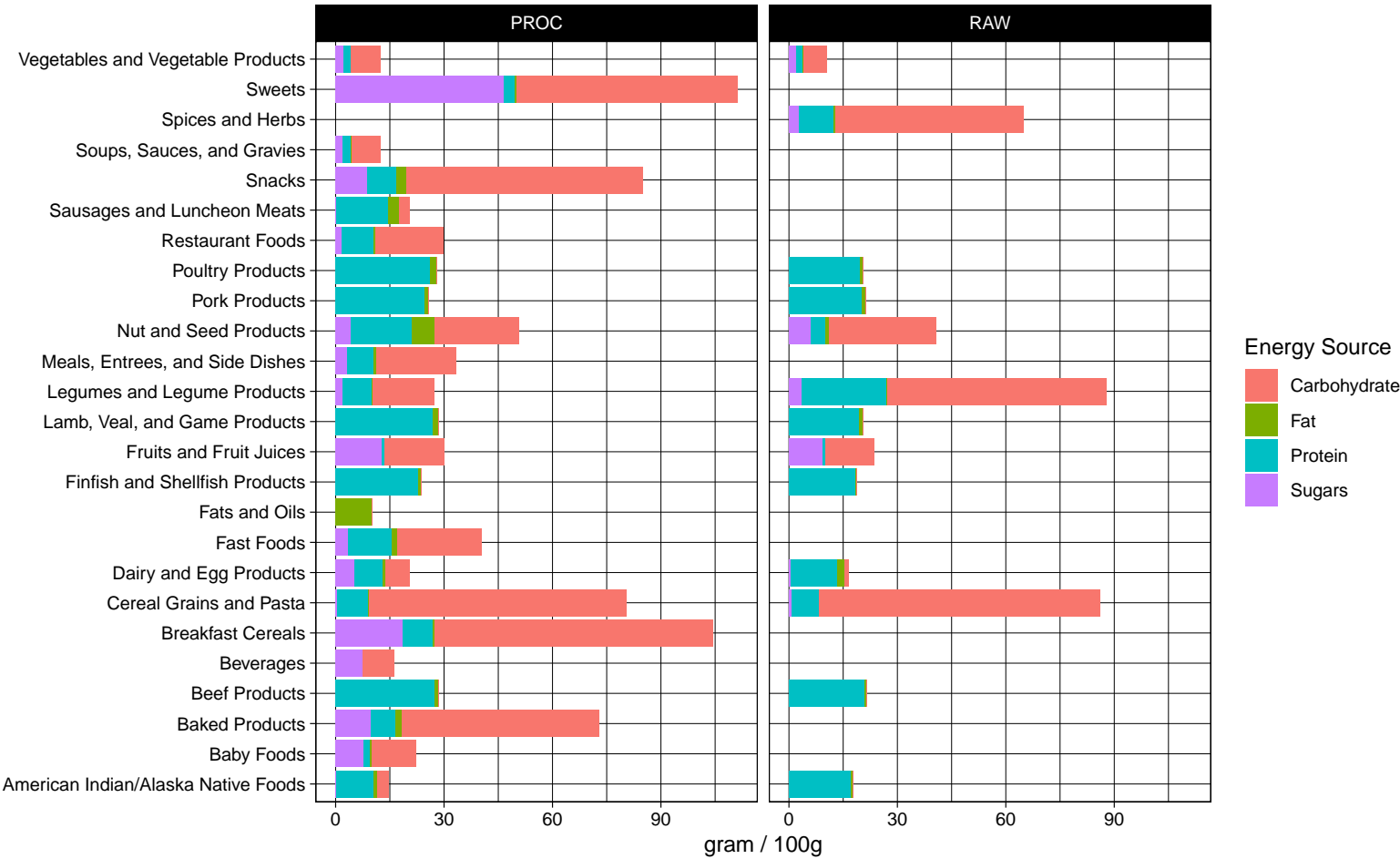


Figure 2: Energy Sources In Food

Fig.5 and Fig.6 show the content of Minerals and Vitamins in each food group and enable us to compare raw and processed foods in that group. scale used is based on the median value for each group. Numerical values shows q75, median and q25 (Upper, meddle and lower).

At first glance we can see that raw foods rarely contains any Fluoride, F and most processed foods contain Fluoride, F. Sodium, Na will range higher in processed products compared to raw foods, but again Legumes and Legume Products values are higher in its raw form compared to its processed form.

Calcium, Ca to my surprise scours the highest in Spices and Herbs in its raw form, with more than double the value compared to Dairy and Egg Products in its processed form. Spices and Herbs have the highest Iron, Fe, Manganese, Mn in all food groups, and high content value for other Minerals, although you cant eat much Spices and Herbs in one day, mixing them with food seems like a nutritious idea.

we can find another difference in Vintamin C as we can see that it rangers higher in raw foods, and its main source are plants sourced food, Vegetables, Fruits, Nuts/Seeds and Spices/Herbs.

above are a small part of our data. we can further explore the ratio of fat and carbohydrates different types in each food group and much more. but we cant explore every detail in this project because this is a much larger subject, and i don't have domain knowledge to do so.

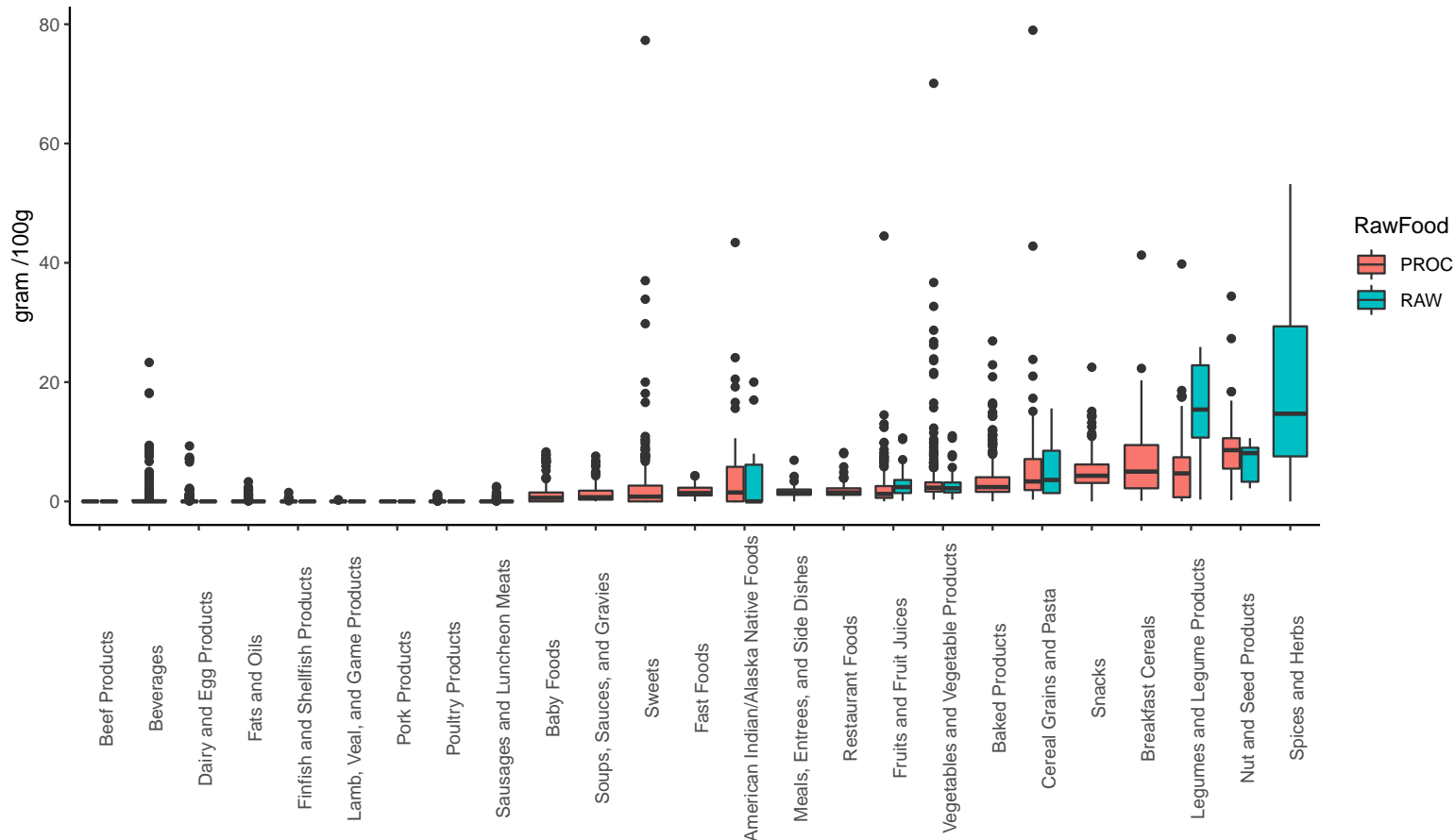


Figure 3: Fiber In Food

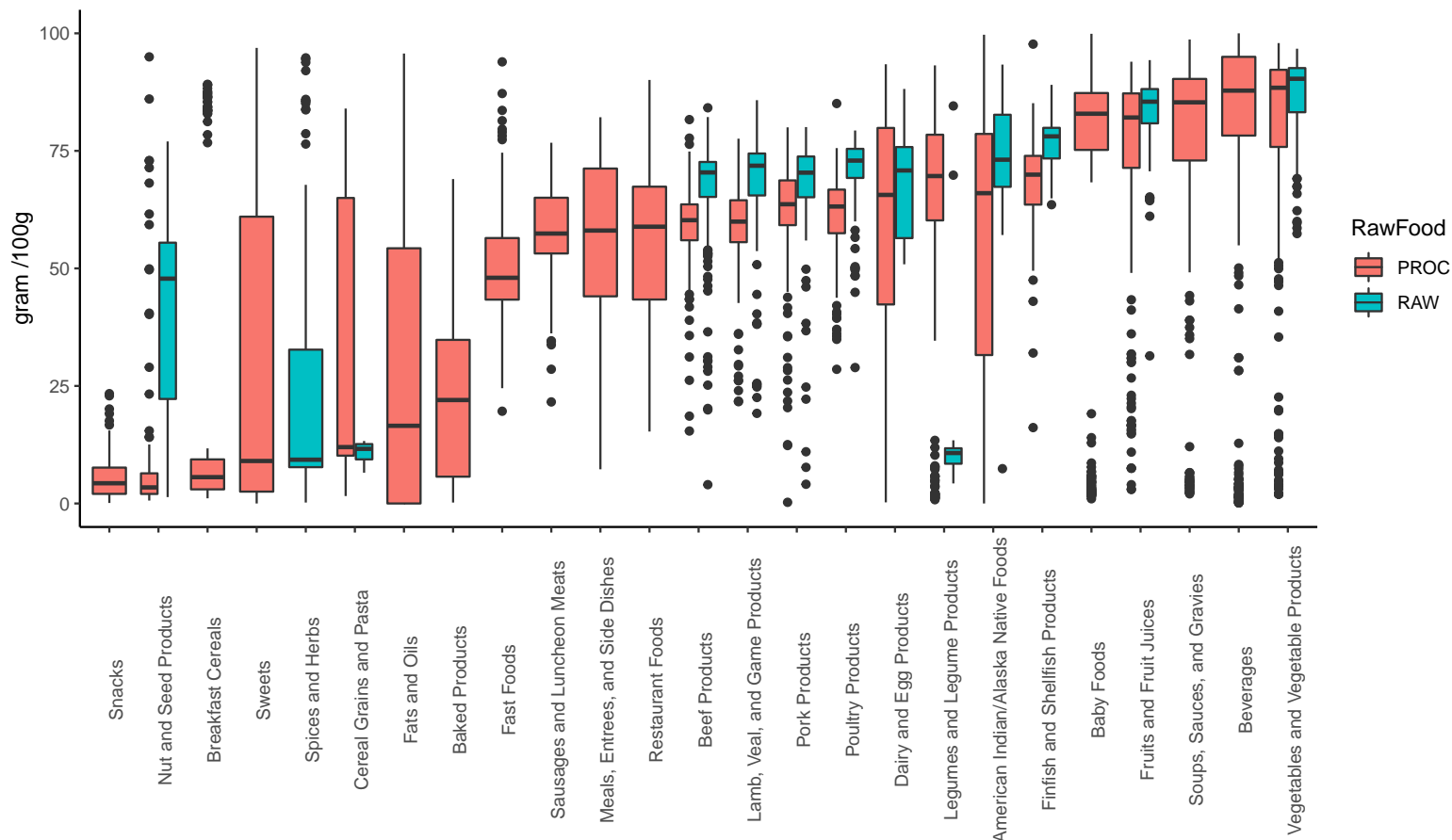


Figure 4: Water In Food

	American Indian/Alaska Native Foods																																																																						
	Baby Foods			Baked Products			Beef Products			Beverages			Breakfast Cereals			Cereal Grains and Pasta			Dairy and Egg Products			Fast Foods			Fats and Oils			Finfish and Shellfish Products			Fruits and Fruit Juices			Lamb, Veal, and Game Products			Legumes and Legume Products			Meals, Entrees, and Side Dishes			Nut and Seed Products			Pork Products			Poultry Products			Restaurant Foods			Sausages and Luncheon Meats			Snacks			Soups, Sauces, and Gravies			Spices and Herbs			Sweets			Vegetables and Vegetable Products	
Calcium, Ca	53.75	125	102.5	16	40.25	242.5	34	574.75	154	13.5	73.5	21	22	123	108.75	141	19	19	89	42	135	33	0	117	48.75																																														
	22	34	52	12	5	55	17	141.5	105	1	37	11	14	54.5	48.5	71	10	14	34	15	58	17	0	54.5	26																																														
	7	11	26	7	2	19	8	94	51	0	14	8	8	35	29	46	7	12	17	10	21	10	0	8.75	14																																														
RAW	24.5	0	0	16	0	0	28.5	113	0	0	47	24	15	174.5	0	54.25	16	13	0	0	0	0	1139.5	0	67																																														
	13	0	0	12	0	0	21	64	0	0	25	12	10	130.5	0	33.5	12	11	0	0	0	0	383	0	35.5																																														
	6.25	0	0	5	0	0	9	56	0	0	12	8	6	89.75	0	17.5	6.5	7.5	0	0	0	0	112.5	0	17																																														
Phosphorus, P	62	21	40	24	11	106	93.75	28	25	2	42.5	15	26	74	30.25	304.25	24	26	28	19	109.25	17	0	44.25	30																																														
	27	10	25	22	3	61	40.5	16	22	0	36	9	24	46	20	217.5	21	23	22	16	75.5	8	0	14	20																																														
	19	6	16	20	1	27	21	11	19	0	30	5	22	29	15.75	109.75	18	20	17	13	50.75	4	0	3	13																																														
RAW	27.5	0	0	23	0	0	129.5	13	0	0	34	17	24	196.5	0	87	23	24	0	0	0	0	285	0	38																																														
	25	0	0	20	0	0	79	11	0	0	30	11	22	183	0	55	21	22	0	0	0	0	174.5	0	23																																														
	21.5	0	0	17	0	0	30	9	0	0	25	8	19	157	0	32	16.5	19	0	0	0	0	73.75	0	14																																														
Potassium, K	303.5	109	202.75	233	31.5	346.5	274.25	477.25	225	16	304.5	25	237	231	165.75	666	267	236.75	235	192	292	65	0	142	70																																														
	170	41	122	210	5	222	119.5	136.5	184	0	257	12	210	142	117.5	456	237	203	130	157	207	35	0	82	42.5																																														
	59.5	20	89	183	1	81	76	98.5	131	0	215.5	9	187.25	99	72	202	216.25	169.5	89	129.5	145	20	0	8	27																																														
RAW	216.5	0	0	213	0	0	298	390	0	0	245	26	208	439.5	0	178.25	227.5	216.5	0	0	0	0	316	0	75.25																																														
	150	0	0	198	0	0	221	208	0	0	210	17	186	396.5	0	117.5	207	184	0	0	0	0	173.5	0	49.5																																														
	40	0	0	175	0	0	111.5	186	0	0	187	12	163	367	0	89.5	186	160	0	0	0	0	69	0	30.75																																														
Sodium, Na	458	200	201.25	365	116.5	358	287	219	250	36.75	469	200	339	438.75	242.75	733	388	266	305	323.25	412.75	255	0	284	378.75																																														
	309	125	136.5	316.5	25	226	140	152	202	10	353	128	306	311	191.5	582	344	242.5	187	252.5	286	140	0	143	249.5																																														
	114	77	103	272.5	8	120.5	56.5	115	177	0	251.5	89	248	155.5	152.25	424	275	199	136	192.75	216.75	71	0	36.5	164																																														
RAW	370	0	0	343.75	0	0	274	163	0	0	394	266	323.25	1405.25	0	569.25	376.5	256	0	0	0	0	1441	0	418.5																																														
	317	0	0	317	0	0	195	135	0	0	350	178.5	288.5	1250	0	465.5	341	228	0	0	0	0	1000	0	317																																														
	148	0	0	283	0	0	100.5	121	0	0	261	126.5	237.75	978.5	0	364.25	280	204	0	0	0	0	458	0	228.5																																														
Fluoride, F	220.5	49	580	69	48.25	537.5	12	628	671.25	779	399	7	84	281.5	591	143	988	237	618	1149.5	538.5	708.25	0	188.75	255																																														
	71	25.5	438	63	13	344	5	112	566	28	124	4	72	65	423	28	100	99	430	960	318	406.5	0	79	120.5																																														
	13.75	6.75	327	54.75	4	13.5	2	53	432.75	0	73.5	2	61	10	349	4	67	75	344	813.5	204.5	306.5	0	33.25	14																																														
RAW	54.5	0	0	67	0	0	5	166	0	0	104	4	83	17.75	0	16.25	87	113.5	0	0	0	0	118	0	27																																														
	43	0	0	58	0	0	5	142	0	0	70	1.5	70.5	13	0	6	63	77	0	0	0	0	50	0	13																																														
	8.5	0	0	53	0	0	2.5	128	0	0	52	1	59	7.25	0	2.75	52.5	63	0	0	0	0	19	0	5																																														
Selenium, Se	0	21.125	49.15	22.4	78.65	71.6	18	12	30.475	27.525	140.25	62.95	26.5	9.7	56.5	7.5	0	14.7	0	41.2	2.5	65.5	0	25.7	37																																														
	0	8.15	25.05	22.4	60.6	53.35	7	8.7	29.65	22.55	31	19.8	21	3.1	56.5	5	0	14.7	0	36	0.1	42.7	0	7	25.6																																														
	0	1.025	15.25	22.4	34.75	21	6	3.05	28.825	13.675	31	7.15	13	2.4	56.5	2.5	0	14.7	0	20.3	0	35.95	0	1.6	6.9																																														
RAW	0	0	0	22.4	0	0	0	1.1	0	0	33.025	4	0	2.2	0	3.4	0	0	0	0	0	0	18.1	0	3.6																																														
	0	0	0	22.4	0	0	0	1.1	0	0	25.55	2.2	0	2.2	0	3.4	0	0	0	0	0	0	2	0	2																																														
	0	0	0	22.4	0	0	0	1.1	0	0	18.075	1.5	0	2.2	0	3.4	0	0	0	0	0	0	1.75	0	1.2																																														
Iron, Fe	53.575	8.4	23.4	33.3	0.5	21.95	31.9	14.75	24.9	1.6	50.025	0.4	26.175	7.15	20.1	34.4	44.7	31.4	18.85	22.825	15.7	4.85	0	3.2	1																																														
	16.55	1.9	13.1	28.7	0.1	12.8	15.5	7.35	20.1	0	45.3	0.3	13.5	2.8	13	11	37.9	26.55	15.2	17.45	8.95	2.8	0	1.9	0.7																																														
	5.25	0.6	6.1	22.875	0	6.15	8.3	2.6	16.575	0	32.75	0.1	9.3	1.3	9.5	4.1	27.5	21.275	6.95	12.5	5.8	1.1	0	0.7	0.4																																														
RAW	30.45	0	0	25.2	0	0	22.25	53.5	0	0	36.5	0.6	19.275	10.925	0	10.1	34.6	22.9	0	0	0	0	12.45	0	0.9																																														
	14.7	0	0	22.2	0	0	15.1	36.4	0	0	36.5	0.3	8.9	8.2	0	7	29.4	20.6	0	0	0	0	5.1	0	0.7																																														
	4.35	0	0	19.2	0	0	12.025	30.7	0	0	23.7	0.1	5.5	4.225	0	6.2	22.65	15.5	0	0	0	0	2.95	0	0.5																																														
Zinc, Zn	4.715	2.38	3.43	3.062	0.25	28.6	4.188	0.42	2.32	0.2	1.455	0.75	2.388	2.59	2.13	5.4	1.13	1.677	1.22	1.768	3.185	0.94	0	1.152	1.268																																														
	2.6	0.8	2.53	2.67	0.06	9.09	2.015	0.16	1.765	0.02	0.85	0.35	2.03	1.7	1.34	3.53	0.94	1.305	0.84	1.24	1.89	0.58	0	0.405	0.73																																														
	0.988	0.3	1.62	2.277	0.02	3.48	1.05	0.05	0.998	0	0.495	0.252	1.518	1.105	0.95	2.53	0.82	0.98	0.66	0.94	1.42	0.288	0	0.09	0.462																																														
RAW	4.7	0	0	2.27	0	0	3.66	3.7	0	0	1.13	0.55	1.918	8.252	0	2.8	0.985	2.61	0	0	0	0	32.125	0	1.7																																														
	1.8	0	0	1.97	0	0	1.96	3.4	0	0	0.59	0.25	1.55	6.6	0	1.865	0.84	1.09	0	0	0	0	13.97	0	0.85																																														
	0.5	0	0	1.69	0	0	1.045	1.74	0	0	0.34	0.15	1.205	5.005	0	1.008	0.615	0.77	0	0	0	0	4.49	0	0.47																																														
Copper, Cu	3.195	1.15	1.01	7.393	0.13	6.755	1.92	2.768	1.725	0.11	1.295	0.18	5.27	1.8	1.058	5.29	2.94	2.982	1.16	2.545	2.64	0.64	0	0.8																																															

[illegible]

3 Building Our Model

We want to build a model that can classify if a food product is **Raw** or **Processed** using its nutrition values. We will use tree based models, because they are easy to use and they will reveal more differences between raw and processed foods by finding variable importance in our model. we will split our data set based on **RawFood** field into training and test sets. any nutrients defined as **added** will be deleted. Below shows the structure of our training set.

we will compare *ctree* from *party* package, *gbm* from *gbm* package and *ranger* from *ranger* package. we will not use **FoodGroup** field to build our model, but we will keep it to compare the performance of our models through different food groups.

```
## Classes 'tbl_df', 'tbl' and 'data.frame':   3895 obs. of  97 variables:
##  $ FoodGroup      : Factor w/ 25 levels "American Indian/Alaska Native Foods",...: 8 8 8 8 8
##  $ RawFood         : Factor w/ 2 levels "PROC","RAW": 1 1 1 1 1 1 1 1 1 1 ...
##  $ Alanine         : num  0.029 0.67 0.659 0.409 0.639 0.823 0.958 0.691 0.741 0.757 ...
##  $ Alcohol, ethyl  : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Arginine        : num  0.031 0.874 0.883 0.467 0.47 ...
##  $ Ash             : num  2.11 3.18 3.6 1.27 5.2 3.79 4.3 3.55 3.27 3.83 ...
##  $ Aspartic acid   : num  0.064 1.588 1.502 0.963 0.779 ...
##  $ Beta-sitosterol : num  4 0 0 0 0 0 0 0 0 0 ...
##  $ Betaine         : num  0.3 0 0 0.6 0 0 0 0 0 0.7 ...
##  $ Caffeine        : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Calcium, Ca     : num  24 674 643 111 493 ...
##  $ Campesterol     : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Carbohydrate, by difference : num  0.06 2.79 4.78 4.76 3.88 1.55 0.36 0.68 2.77 5.58 ...
##  $ Carotene, alpha : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Carotene, beta  : num  158 76 0 6 3 32 33 78 41 51 ...
##  $ Cholesterol      : num  215 94 103 12 89 116 110 89 64 64 ...
##  $ Choline, total   : num  18.8 15.4 0 16.3 15.4 15.4 15.4 15.4 14.2 ...
##  $ Copper, Cu       : num  0 0.024 0.042 0.033 0.032 0.025 0.032 0.032 0.025 0.033 ...
##  $ Cryptoxanthin, beta : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Cystine          : num  0.008 0.131 0.117 0.062 0.083 0.261 0.304 0.123 0.144 0.124 ...
##  $ Dihydrophyllquinone : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Energy           : num  717 371 387 81 265 389 413 373 254 295 ...
##  $ Fatty acids, total monounsaturated: num  21.021 8.598 8.671 0.516 4.623 ...
##  $ Fatty acids, total polyunsaturated: num  3.043 0.784 0.87 0.083 0.591 ...
##  $ Fatty acids, total saturated      : num  51.37 18.76 19.48 1.24 13.3 ...
##  $ Fatty acids, total trans           : num  3.278 0 0 0.067 0 ...
##  $ Fatty acids, total trans-monoenoic: num  2.982 0 0 0.054 0 ...
##  $ Fatty acids, total trans-polyenoic: num  0.296 0 0 0.013 0 0 0 0 0.138 ...
##  $ Fiber, total dietary              : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Fluoride, F                      : num  2.8 0 0 0 0 0 0 0 0 0 ...
##  $ Folate, DFE                     : num  3 20 18 8 32 6 10 18 9 27 ...
##  $ Folate, food                    : num  3 20 18 8 32 6 10 18 9 27 ...
##  $ Folate, total                   : num  3 20 18 8 32 6 10 18 9 27 ...
##  $ Folic acid                     : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Fructose                       : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Galactose                      : num  0 0 0 0.12 0 0 0 0 0 0.78 ...
##  $ Glucose (dextrose)             : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Glutamic acid                  : num  0.178 5.515 5.718 2.446 2.421 ...
##  $ Glycine                       : num  0.018 0.437 0.403 0.209 0.097 0.457 0.533 0.422 0.464 0.551 ...
##  $ Histidine                     : num  0.023 0.823 0.821 0.306 0.397 ...
##  $ Hydroxyproline                 : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Iron, Fe                      : num  0.02 0.43 0.21 0.13 0.65 0.23 0.17 0.72 0.22 0.22 ...
##  $ Isoleucine                    : num  0.051 1.137 1.451 0.556 0.803 ...
##  $ Lactose                       : num  0 0 0 3.87 0 0 0 0 0 1.12 ...
##  $ Leucine                       : num  0.083 2.244 2.238 1.049 1.395 ...
##  $ Lutein + zeaxanthin           : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Lycopene                      : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Lysine                        : num  0.067 2.124 1.945 0.878 1.219 ...
##  $ Magnesium, Mg                 : num  2 24 21 9 19 14 36 27 23 27 ...
##  $ Maltose                       : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ Manganese, Mn                 : num  0 0.012 0.012 0.015 0.028 0.014 0.017 0.011 0.01 0.039 ...
##  $ Menaquinone-4                 : num  0 0 0 0 0 0 0 0 0 4.1 ...
##  $ Methionine                    : num  0.021 0.565 0.612 0.253 0.368 0.706 0.822 0.641 0.677 0.551 ...
```

## \$ Niacin	: num	0.042	0.118	0.08	0.103	0.991	0.15	0.106	0.093	0.105	0.111	...
## \$ Pantothenic acid	: num	0.11	0.288	0.413	0.524	0.967	0.429	0.562	0.21	0.079	0.429	...
## \$ Phenylalanine	: num	0.041	1.231	1.231	0.543	0.675	...					
## \$ Phosphorus, P	: num	24	451	464	150	337	346	605	444	463	548	...
## \$ Phytosterols	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Potassium, K	: num	24	136	95	125	62	64	81	81	84	188	...
## \$ Proline	: num	0.082	2.575	2.634	1.155	1.378	...					
## \$ Protein	: num	0.85	23.24	23.37	10.45	14.21	...					
## \$ Retinol	: num	671	286	220	68	125	258	268	192	124	223	...
## \$ Riboflavin	: num	0.034	0.351	0.293	0.251	0.844	0.204	0.279	0.39	0.303	0.353	...
## \$ Selenium, Se	: num	1	14.5	14.5	11.9	15	14.5	14.5	14.5	14.4	27.6	...
## \$ Serine	: num	0.046	1.289	1.366	0.601	1.169	...					
## \$ Sodium, Na	: num	643	560	700	308	1139	...					
## \$ Starch	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Stigmasterol	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Sucrose	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Sugars, total	: num	0.06	0.51	0	4	0	1.55	0.36	0.5	1.13	1.9	...
## \$ Theobromine	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Thiamin	: num	0.005	0.014	0.046	0.02	0.154	0.021	0.06	0.015	0.018	0.024	...
## \$ Threonine	: num	0.038	0.882	0.832	0.47	0.637	...					
## \$ Tocopherol, beta	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Tocopherol, delta	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Tocopherol, gamma	: num	0	0	0	0	0	0	0	0	0.03	...	
## \$ Tocotrienol, alpha	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Tocotrienol, beta	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Tocotrienol, delta	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Tocotrienol, gamma	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Total lipid (fat)	: num	81.11	29.68	30.6	2.27	21.49	...					
## \$ Tryptophan	: num	0.012	0.324	0.3	0.138	0.2	0.361	0.421	0.315	0.339	0.551	...
## \$ Tyrosine	: num	0.041	1.115	1.128	0.568	0.668	...					
## \$ Valine	: num	0.057	1.472	1.56	0.703	1.065	...					
## \$ Vitamin A, IU	: num	2499	1080	985	236	422	...					
## \$ Vitamin A, RAE	: num	684	292	233	68	125	261	271	198	127	227	...
## \$ Vitamin B-12	: num	0.17	1.26	0.83	0.47	1.69	1.68	1.6	0.83	0.82	1.68	...
## \$ Vitamin B-6	: num	0.003	0.065	0.074	0.057	0.424	0.083	0.081	0.079	0.07	0.1	...
## \$ Vitamin C, total ascorbic acid	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Vitamin D	: num	0	22	0	0	16	23	24	22	12	15	...
## \$ Vitamin D (D2 + D3)	: num	0	0.5	0	0	0.4	0.6	0.6	0.6	0.3	0.4	...
## \$ Vitamin D2 (ergocalciferol)	: num	0	0	0	0	0	0	0	0	0	0	...
## \$ Vitamin D3 (cholecalciferol)	: num	0	0.5	0	0	0.4	0.6	0.6	0.6	0.3	0.4	...
## \$ Vitamin E (alpha-tocopherol)	: num	2.32	0.26	0	0.08	0.18	0.27	0.28	0.26	0.14	0.5	...
## \$ Vitamin K (phylloquinone)	: num	7	2.5	0	0	1.8	2.6	2.7	2.5	1.6	1.3	...
## \$ Water	: num	16.2	41.1	37.6	81.2	55.2	...					
## \$ Zinc, Zn	: num	0.09	2.6	2.79	0.51	2.88	3.5	3.9	3	2.76	3.62	...

Table 6: Model Performance, Training Set

	ctree_Training	gbm_Training	ranger_Training	ctree_Test	gbm_Test	ranger_Test
Sensitivity	0.811	0.842	1.000	0.742	0.758	0.770
Specificity	0.968	0.978	1.000	0.947	0.970	0.975
Pos Pred Value	0.848	0.896	1.000	0.758	0.850	0.872
Neg Pred Value	0.958	0.965	1.000	0.943	0.948	0.950
Precision	0.848	0.896	1.000	0.758	0.850	0.872
Recall	0.811	0.842	1.000	0.742	0.758	0.770
F1	0.829	0.868	1.000	0.749	0.801	0.818
Prevalence	0.182	0.182	0.182	0.182	0.182	0.182
Detection Rate	0.147	0.153	0.182	0.135	0.138	0.140
Detection Prevalence	0.174	0.171	0.182	0.178	0.162	0.160
Balanced Accuracy	0.889	0.910	1.000	0.844	0.864	0.872

Table.6 shows the performance of each model, we can see that **ranger** has the highest accuracy of **1**, followed by **gbm**. a perfect accuracy may be a result of over fitting. applying our models to the test set we can see that accuracy of **ranger** falls to **.87** and **gbm** falls to **.86**. which indicates that **gbm** model is more consistent and not over fitted.

One more thing to notice that if we set **RAW** as our positive test result we get a **Prevalence** of **0.182**, thus the high **Specificity** relative to **Sensitivity**. trying to balance our data set by sampling up and down we get the results in *table.7*.

Table 7: Model Performance, Balanced Data Set

	up	gbm	down
Sensitivity	0.968	0.842	0.980
Specificity	0.930	0.978	0.901
Pos Pred Value	0.754	0.896	0.687
Neg Pred Value	0.992	0.965	0.995
Precision	0.754	0.896	0.687
Recall	0.968	0.842	0.980
F1	0.847	0.868	0.808
Prevalence	0.182	0.182	0.182
Detection Rate	0.176	0.153	0.178
Detection Prevalence	0.233	0.171	0.259
Balanced Accuracy	0.949	0.910	0.941

as we can see in *Table.8*, sampling up/down increased **Balanced Accuracy** and increased **Sensitivity** on the expense of **Specificity**, **F1** and **Precision**. this means that our model started to predict more **raw foods**, enough to correctly cover *0.97-0.98* of the **RAW** population but on the expense of predicting more **processed foods** as **raw foods**.

Below shows a breakdown for the performance of **gbm**, and indeed we see that sampling enhanced **Sensitivity**, but decreased **Specificity**. so in conclusion **gbm** is the more stable model and to chose to sample up, down or not to sample at all depends on the importance of predicting positive(raw) test results more than negative results. but in our case i would prefer a balanced model and will chose **gbm** without sampling.

Table 8: GBM, Model Performance Per Group

FoodGroup	Sensitivity	Specificity	Precision	Recall	F1	Prevalence	Balanced Accuracy
Dairy and Egg Products	0.71	1.00	1.00	0.71	0.83	0.05	0.86
Spices and Herbs	0.77	NA	1.00	0.77	0.87	1.00	NA
Baby Foods	NA	0.96	0.00	NA	NA	0.00	NA
Fats and Oils	NA	1.00	NA	NA	NA	0.00	NA
Poultry Products	0.92	0.98	0.96	0.92	0.94	0.32	0.95
Soups, Sauces, and Gravies	NA	1.00	NA	NA	NA	0.00	NA
Sausages and Luncheon Meats	NA	0.99	0.00	NA	NA	0.00	NA
Breakfast Cereals	NA	1.00	NA	NA	NA	0.00	NA
Snacks	NA	1.00	NA	NA	NA	0.00	NA
Fruits and Fruit Juices	0.67	0.99	0.98	0.67	0.80	0.34	0.83
Pork Products	0.81	0.99	0.97	0.81	0.89	0.26	0.90
Vegetables and Vegetable Products	0.67	0.93	0.77	0.67	0.72	0.25	0.80
Nut and Seed Products	0.33	1.00	1.00	0.33	0.50	0.14	0.67
Beef Products	0.99	0.94	0.92	0.99	0.95	0.41	0.96
Beverages	NA	1.00	NA	NA	NA	0.00	NA
Finfish and Shellfish Products	0.88	0.88	0.86	0.88	0.87	0.45	0.88
Legumes and Legume Products	0.92	0.99	0.96	0.92	0.94	0.16	0.95
Lamb, Veal, and Game Products	0.98	0.94	0.93	0.98	0.95	0.43	0.96
Baked Products	NA	1.00	NA	NA	NA	0.00	NA
Sweets	NA	1.00	NA	NA	NA	0.00	NA
Cereal Grains and Pasta	0.00	1.00	NA	0.00	NA	0.04	0.50
Fast Foods	NA	1.00	NA	NA	NA	0.00	NA
Meals, Entrees, and Side Dishes	NA	1.00	NA	NA	NA	0.00	NA
American Indian/Alaska Native Foods	0.38	0.98	0.89	0.38	0.53	0.24	0.68
Restaurant Foods	NA	1.00	NA	NA	NA	0.00	NA

Table 9: GBM, Model Performance Per Group, Sampling UP

FoodGroup	Sensitivity	Specificity	Precision	Recall	F1	Prevalence	Balanced Accuracy
Dairy and Egg Products	0.86	1.00	1.00	0.86	0.92	0.05	0.93
Spices and Herbs	0.91	NA	1.00	0.91	0.95	1.00	NA
Baby Foods	NA	0.94	0.00	NA	NA	0.00	NA
Fats and Oils	NA	0.99	0.00	NA	NA	0.00	NA
Poultry Products	1.00	0.91	0.85	1.00	0.92	0.32	0.96
Soups, Sauces, and Gravies	NA	0.99	0.00	NA	NA	0.00	NA
Sausages and Luncheon Meats	NA	0.99	0.00	NA	NA	0.00	NA
Breakfast Cereals	NA	0.99	0.00	NA	NA	0.00	NA
Snacks	NA	1.00	NA	NA	NA	0.00	NA
Fruits and Fruit Juices	0.95	0.91	0.85	0.95	0.90	0.34	0.93
Pork Products	1.00	0.94	0.84	1.00	0.91	0.26	0.97
Vegetables and Vegetable Products	0.97	0.76	0.58	0.97	0.72	0.25	0.87
Nut and Seed Products	1.00	0.91	0.64	1.00	0.78	0.14	0.95
Beef Products	0.99	0.90	0.87	0.99	0.93	0.41	0.94
Beverages	NA	0.99	0.00	NA	NA	0.00	NA
Finfish and Shellfish Products	0.94	0.75	0.76	0.94	0.84	0.45	0.85
Legumes and Legume Products	1.00	0.90	0.67	1.00	0.80	0.16	0.95
Lamb, Veal, and Game Products	0.99	0.85	0.83	0.99	0.90	0.43	0.92
Baked Products	NA	1.00	0.00	NA	NA	0.00	NA
Sweets	NA	1.00	NA	NA	NA	0.00	NA
Cereal Grains and Pasta	0.00	0.96	0.00	0.00	NaN	0.04	0.48
Fast Foods	NA	1.00	NA	NA	NA	0.00	NA
Meals, Entrees, and Side Dishes	NA	1.00	NA	NA	NA	0.00	NA
American Indian/Alaska Native Foods	0.81	0.80	0.57	0.81	0.67	0.24	0.80
Restaurant Foods	NA	1.00	NA	NA	NA	0.00	NA

4 Conclusion

Data exploration shows that there is a clear difference between raw and processed food if we consider food as groups, such as the differences in Fluoride, F, Sodium, Na and Vitamin C values between raw and processed food. Our model shows that we can predict raw/processed food with a Balanced Accuracy of **0.91**, *Fig.7* shows the most important differences between raw and processed foods, again we can see that Sodium, Na and Vitamin C are important differences between raw and processed foods.

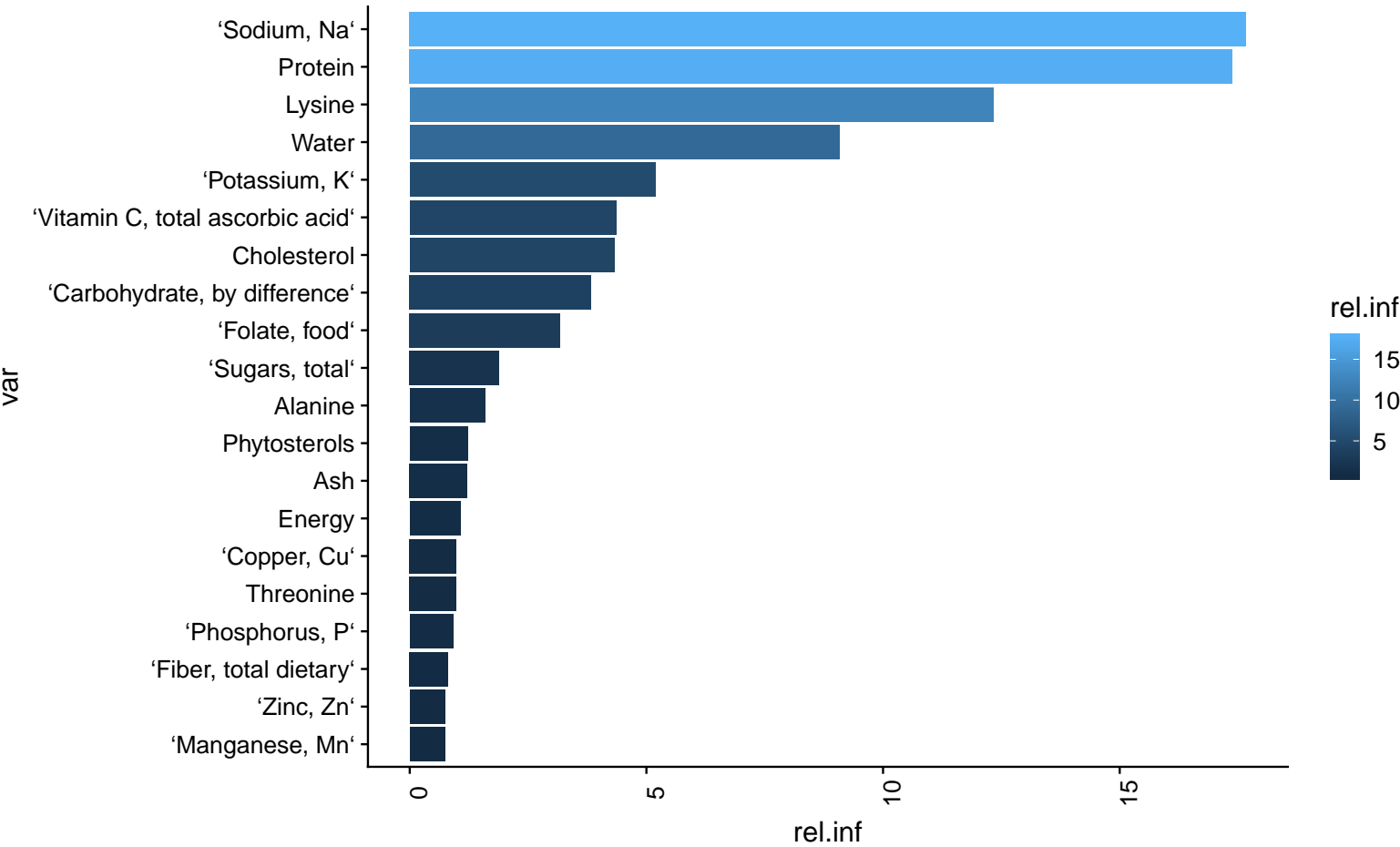


Figure 7: Var Importance