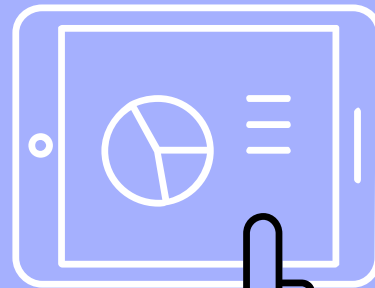
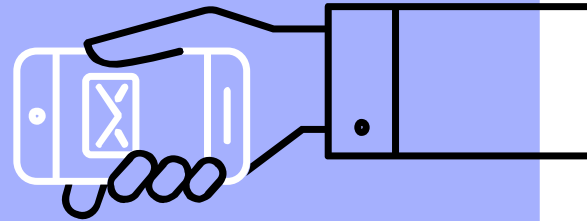
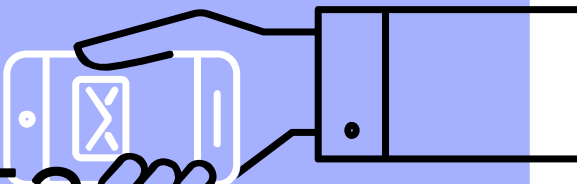
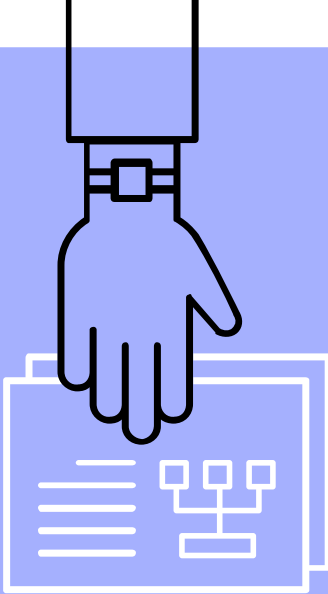




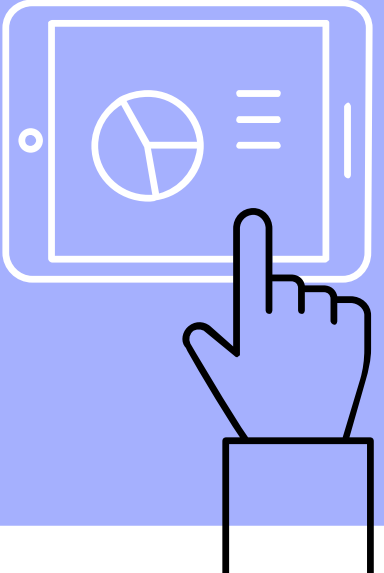
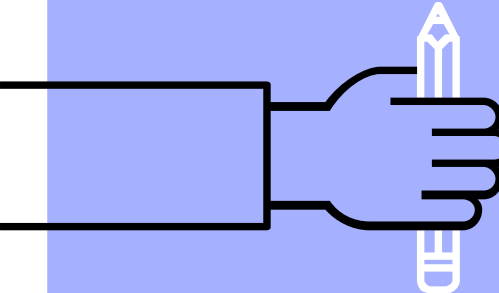
# DATA SCIENCE Olympics

By Smera Gora and ShwinyG on GitHub





# We got our data on Kaggle 270,000 x 15



Variety of different data with categorical and numerical  
predictors

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

freedkn.R olymptics olympticsOmit KNN.R knnOlympics pred1 olympticsKnnTestMedal olympticsMedal

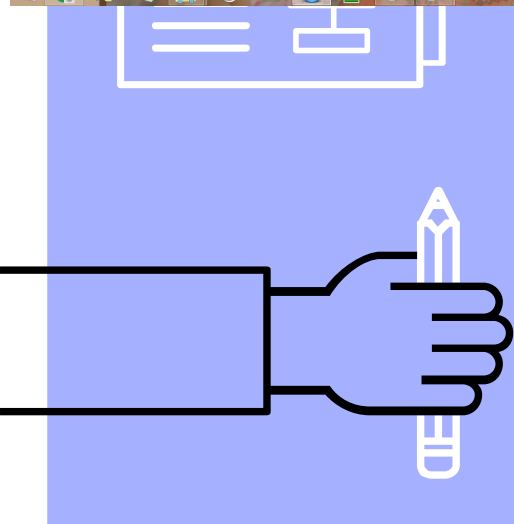
#	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal
1	A Djang	M	24	180	80.0	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	Basketball Men's Basketball	NA
2	A Lamusi	M	23	170	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	NA
3	Gunnar Nielsen Aaby	M	24	NA	NA	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	NA
4	Edgar Lindenau Aabye	M	34	NA	NA	Denmark/Sweden	DEN	1900 Summer	1900	Summer	Paris	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold
5	Christine Jacobs Aaftink	F	21	185	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	NA
6	Christine Jacobs Aaftink	F	21	185	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 1,000 metres	NA
7	Christine Jacobs Aaftink	F	25	185	82.0	Netherlands	NED	1992 Winter	1992	Winter	Albertville	Speed Skating	Speed Skating Women's 500 metres	NA
8	Christine Jacobs Aaftink	F	25	185	82.0	Netherlands	NED	1992 Winter	1992	Winter	Albertville	Speed Skating	Speed Skating Women's 1,000 metres	NA
9	Christine Jacobs Aaftink	F	27	185	82.0	Netherlands	NED	1994 Winter	1994	Winter	Lillehammer	Speed Skating	Speed Skating Women's 500 metres	NA
10	Christine Jacobs Aaftink	F	27	185	82.0	Netherlands	NED	1994 Winter	1994	Winter	Lillehammer	Speed Skating	Speed Skating Women's 1,000 metres	NA
11	Per Knut Aaland	M	31	188	75.0	United States	USA	1992 Winter	1992	Winter	Albertville	Cross Country Skiing	Cross Country Skiing Men's 10 kilometres	NA
12	Per Knut Aaland	M	31	188	75.0	United States	USA	1992 Winter	1992	Winter	Albertville	Cross Country Skiing	Cross Country Skiing Men's 20 kilometres	NA
13	Per Knut Aaland	M	31	188	75.0	United States	USA	1992 Winter	1992	Winter	Albertville	Cross Country Skiing	Cross Country Skiing Men's 10/15 kilometres Pursuit	NA
14	Per Knut Aaland	M	31	188	75.0	United States	USA	1992 Winter	1992	Winter	Albertville	Cross Country Skiing	Cross Country Skiing Men's 4 x 10 kilometres Relay	NA
15	Per Knut Aaland	M	33	188	75.0	United States	USA	1994 Winter	1994	Winter	Lillehammer	Cross Country Skiing	Cross Country Skiing Men's 10 kilometres	NA

Showing 1 to 16 of 271,116 entries

Console Terminal

```
> dplyr::filter(olympics, subswimmer)
> table(dplyr::filter(olympics, subswimmer))
0 1
30188 30181
> #read csv and remove unnecessary rows and getting NAs
> olympics<-read.csv("C:\\Users\\Smera\\Documents\\DTK\\olympics\\Finalproj.R\\olympics.csv", stringsAsFactors = FALSE)
> view(olympics)
>
```

9:38 AM 8/9/2018



RStudio

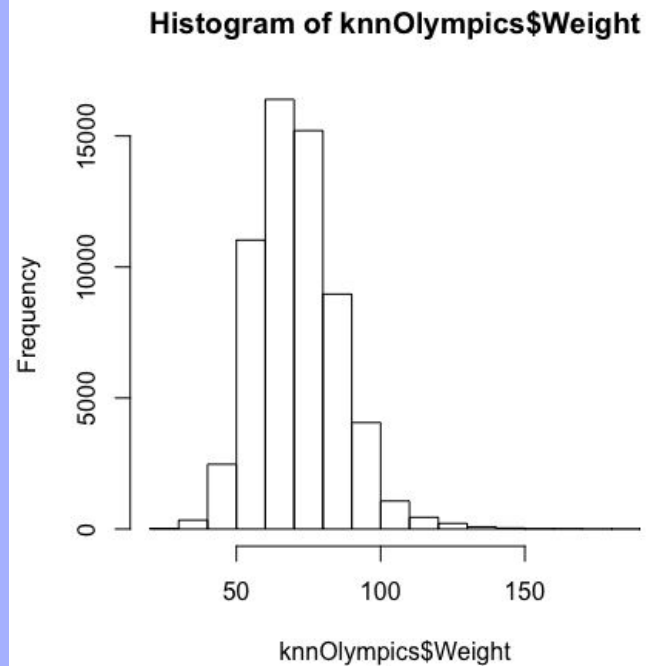
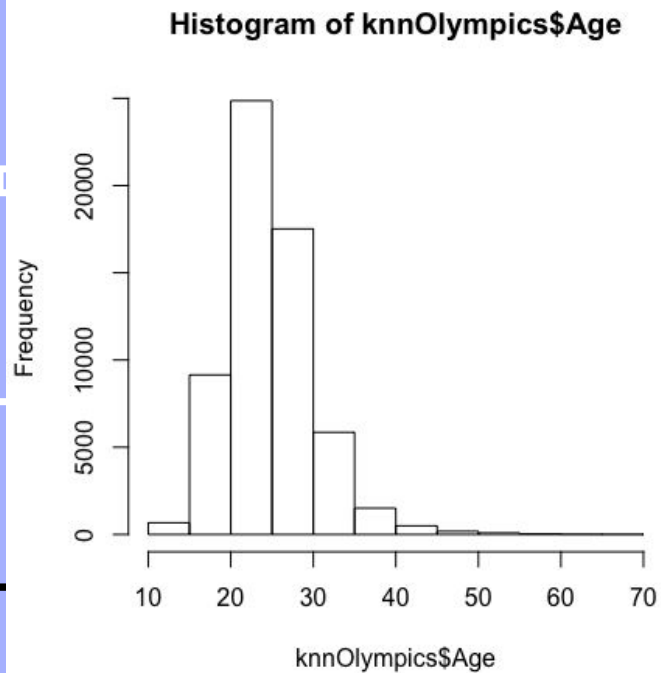
File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

freedkn.R olympticsOmit KNN.R knnOlympics pred1 olympticsKnnTestMedal olympticsMedal

#	Name	Sex	Age	Height	Weight	NOC	Year	Season	Sport	Event	Medal
1	A Djang	M	24	180	80.0	CHN	1992	Summer	Basketball	Basketball Men's Basketball	NA
2	A Lamusi	M	23	170	60.0	CHN	2012	Summer	Judo	Judo Men's Extra-Lightweight	NA
3	Gunnar Nielsen Aaby	M	24	NA	NA	DEN	1920	Summer	Football	Football Men's Football	NA
4	Edgar Lindenau Aabye	M	34	NA	NA	DEN	1900	Summer	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold
5	Christine Jacobs Aaftink	F	21	185	82.0	NED	1988	Winter	Speed Skating	Speed Skating Women's 500 metres	NA
6	Christine Jacobs Aaftink	F	21	185	82.0	NED	1988	Winter	Speed Skating	Speed Skating Women's 1,000 metres	NA
7	Christine Jacobs Aaftink	F	25	185	82.0	NED	1992	Winter	Speed Skating	Speed Skating Women's 500 metres	NA
8	Christine Jacobs Aaftink	F	25	185	82.0	NED	1992	Winter	Speed Skating	Speed Skating Women's 1,000 metres	NA
9	Christine Jacobs Aaftink	F	27	185	82.0	NED	1994	Winter	Speed Skating	Speed Skating Women's 500 metres	NA
10	Christine Jacobs Aaftink	F	27	185	82.0	NED	1994	Winter	Speed Skating	Speed Skating Women's 1,000 metres	NA
11	Per Knut Aaland	M	31	188	75.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 10 kilometres	NA
12	Per Knut Aaland	M	31	188	75.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 20 kilometres	NA
13	Per Knut Aaland	M	31	188	75.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 10/15 kilometres Pursuit	NA
14	Per Knut Aaland	M	31	188	75.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 4 x 10 kilometres Relay	NA
15	Per Knut Aaland	M	33	188	75.0	USA	1994	Winter	Cross Country Skiing	Cross Country Skiing Men's 10 kilometres	NA
16	Per Knut Aaland	M	33	188	75.0	USA	1994	Winter	Cross Country Skiing	Cross Country Skiing Men's 30 kilometres	NA
17	Per Knut Aaland	M	33	188	75.0	USA	1994	Winter	Cross Country Skiing	Cross Country Skiing Men's 10/15 kilometres Pursuit	NA
18	Per Knut Aaland	M	33	188	75.0	USA	1994	Winter	Cross Country Skiing	Cross Country Skiing Men's 4 x 10 kilometres Relay	NA
19	John Aalberg	M	31	183	72.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 10 kilometres	NA
20	John Aalberg	M	31	183	72.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 50 kilometres	NA
21	John Aalberg	M	31	183	72.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 10/15 kilometres Pursuit	NA
22	John Aalberg	M	31	183	72.0	USA	1992	Winter	Cross Country Skiing	Cross Country Skiing Men's 4 x 10 kilometres Relay	NA
23	John Aalberg	M	33	183	72.0	USA	1994	Winter	Cross Country Skiing	Cross Country Skiing Men's 10 kilometres	NA

Showing 1 to 24 of 271,116 entries



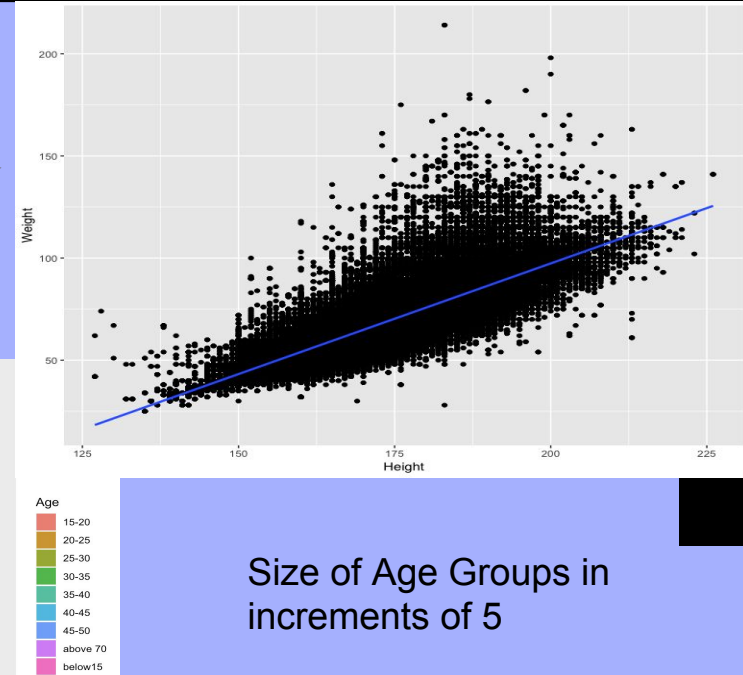
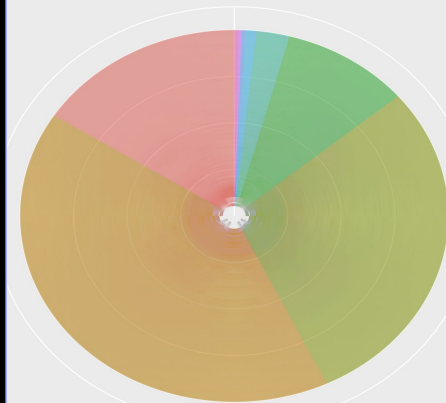
## PIE CHART

- ▶ For KNN and Naive Bayes we had to make all our predictors categorical by splitting them up into groups

## SCATTER PLOT

- ▶ We were able to reduce the dimensions of our df by noticing that height and weight were

Height to Weight  
Correlation=0.796



Size of Age Groups in  
increments of 5

# The three different models

## GLM:

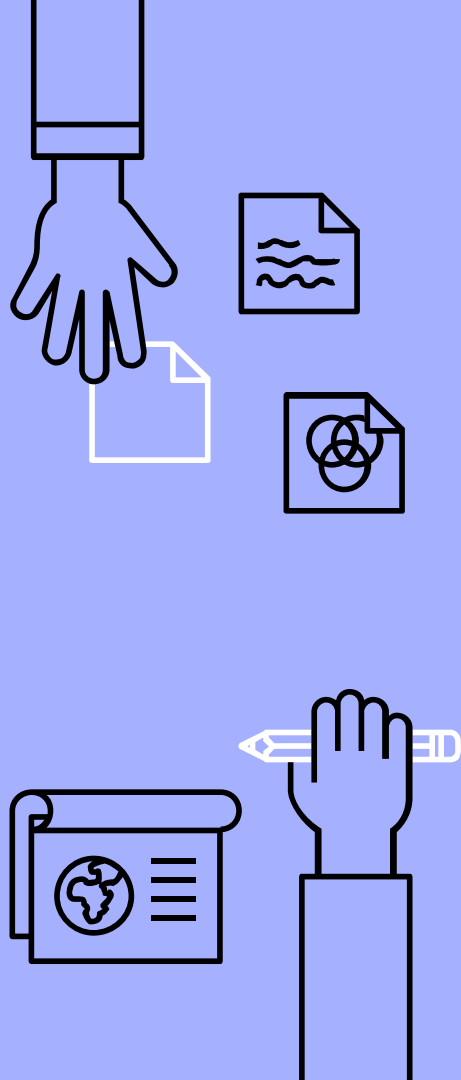
- ▶ Fits **Generalized Linear Models**.
- ▶ Puts coefficients to each predictor
  - A **negative coefficient**: less likely to win
  - A **positive coefficient**: more likely to win

## KNN:

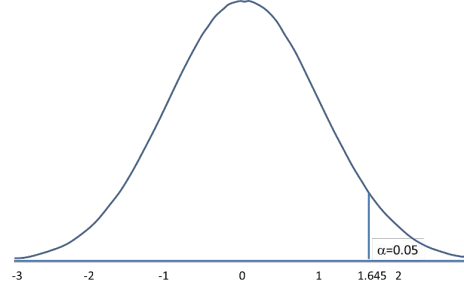
- ▶ **K Nearest Neighbors**
- ▶ K value is determined with seeing the clusters with the **lowest error**
- ▶ **Naive classifier**: belong to the same class
- ▶ **Euclidean Dist**: takes the distance between the points

## Naive Bayes:

- ▶ Computes **probabilities** based on the given **categorical** predictor variables.
- ▶ **Bayes Theorem**

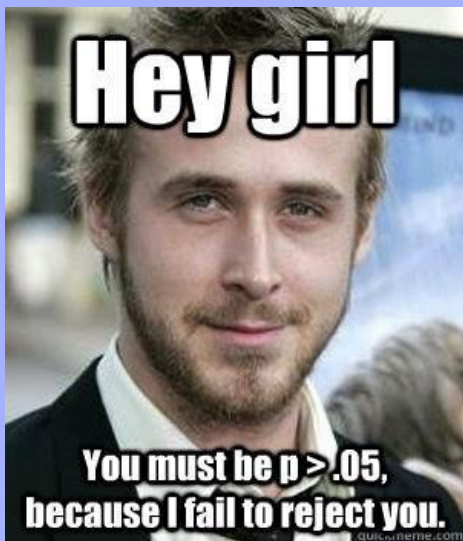


# Vocabulary

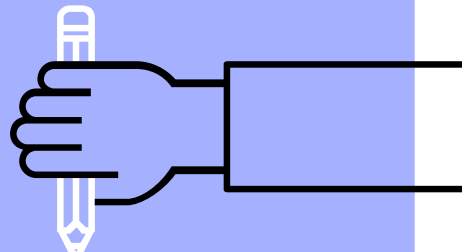


- ▷ **P-Value**: if P is low reject  $H_0$
- ▷ receiver operating characteristic curve, **ROC**: sensitivity vs 1-specificity
- ▷ Area Under Curve, **AUC**: Area Under the ROC. The range is from .5 to 1 and the closer the area is to one the better
- ▷ **True Positive**: predicting the condition positive correctly
- ▷ **True Negative**: predicting condition in negative correctly
- ▷ **False Positive**: Predicting condition is positive incorrectly
- ▷ **False Negative**: Predicting the condition is negative incorrectly
- ▷ **Train**: data set we used to run the model
- ▷ **Test**: data set tested the model on
- ▷ **Accuracy**: correct identification
- ▷ **Sensitivity**: correctly identify true positive
- ▷ **Specificity**: correctly identify true negative
- ▷ **F1**: Measures rate of performance using recall and precision

	Decision	
	Accept $H_0$	Reject $H_0$
$H_0$ (true)	Correct decision	Type I error ( $\alpha$ error)
$H_0$ (false)	Type II error ( $\beta$ error)	Correct decision



# 1. Generalized Linear Model

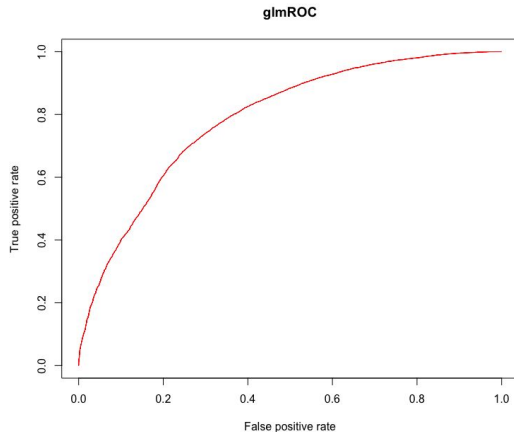




# GLM: Generalized Linear Model

Accuracy: 0.7196  
Sensitivity: 0.6995  
Specificity: 0.6818  
F1: 0.6905

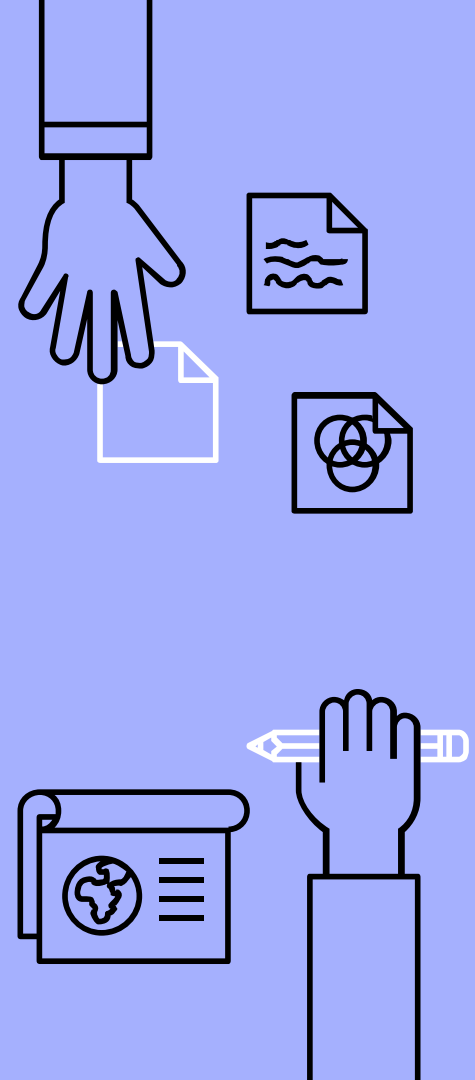
AUC: 0.7201



	False	True
0	6245	2915
1	2162	6785

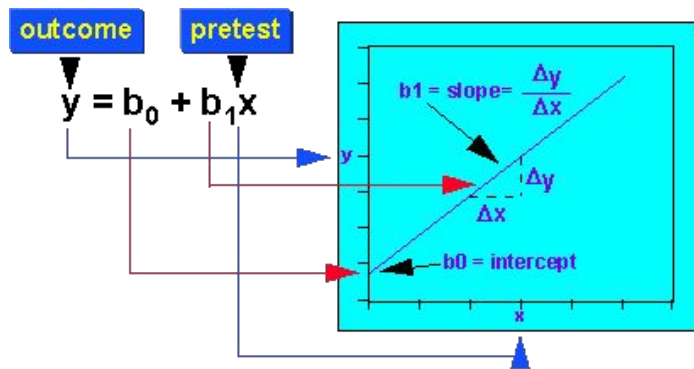
Predictor examples::

- Height
- Gender
- Age
- Women's Archery team
- Women's 4 x 100
- Women's 4 x 400
- Event Archery Men's team
- Men's 1600m relay
- Men's 3 mile
- Men's 3000 miles
- Men's 4 x 100
- Men's 4x400
- Men's CC
- Men's High Jump
- Etc

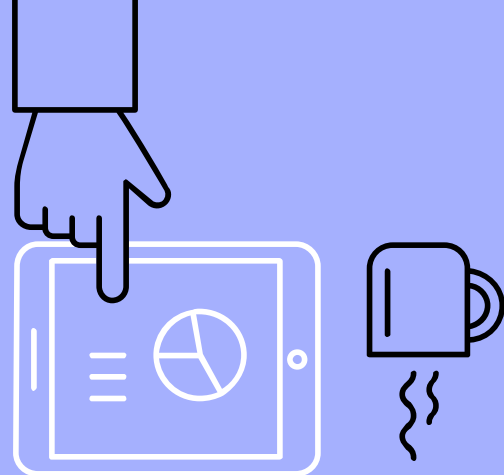


# GLM (problem) cont.

- ▷ Made our response binary
- ▷ Got rid of NAs
- ▷ Ran glm
- ▷ It wouldn't work because we had too many levels in some of our obs variables
- ▷ Once we did some dimension reduction it worked and we got rid of the risk of overfitting
- ▷ We were getting really low sensitivity: it was about 8%
- ▷ The frequency of wins was extremely low compared to the loss counts.
- ▷ Another dataframe with equal frequencies of each



**When you delete a block  
of code that you thought  
was useless**

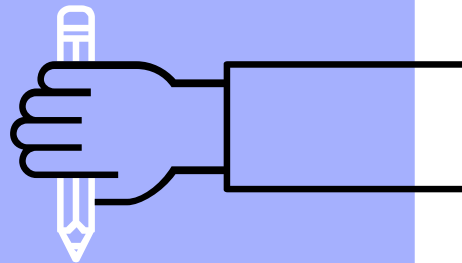


**LOOK AT OLD CODE FROM A  
YEAR AGO TO REFRESH  
MEMORY**

**DIDN'T WRITE A  
SINGLE COMMENT**

**THE ONLY PROGRAMMING  
JOKE I KNOW IS**

**YOUR CODE**



# Naive Bayes

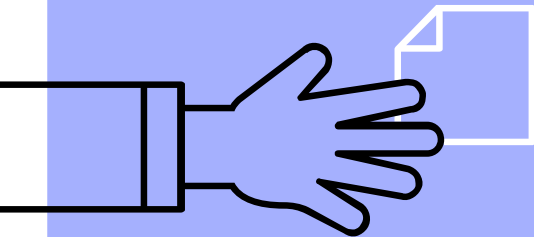
**WROTE R CODE  
NOT USING GOOGLE**



**WORKED ON THE FIRST TRY**

**I AM PROGRAMMER**

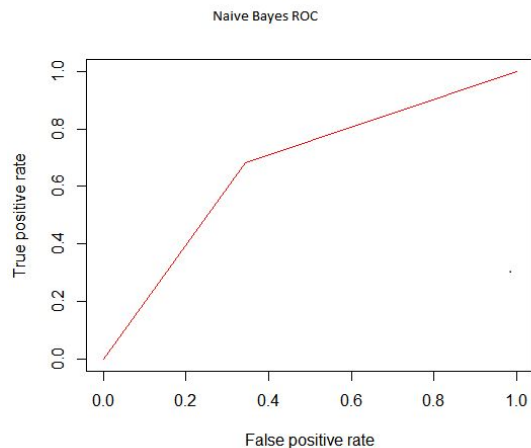
**I MAKE COMPUTER BEEP BOOP  
BEEP BEEP BOOP**



# Naive Bayes

	False	True
0	5981	2837
1	3054	6234

Accuracy: 0.6746  
Sensitivity: 0.6872451  
Specificity: 0.6782717  
F1: 0.6827289  
AUC: 0.6692



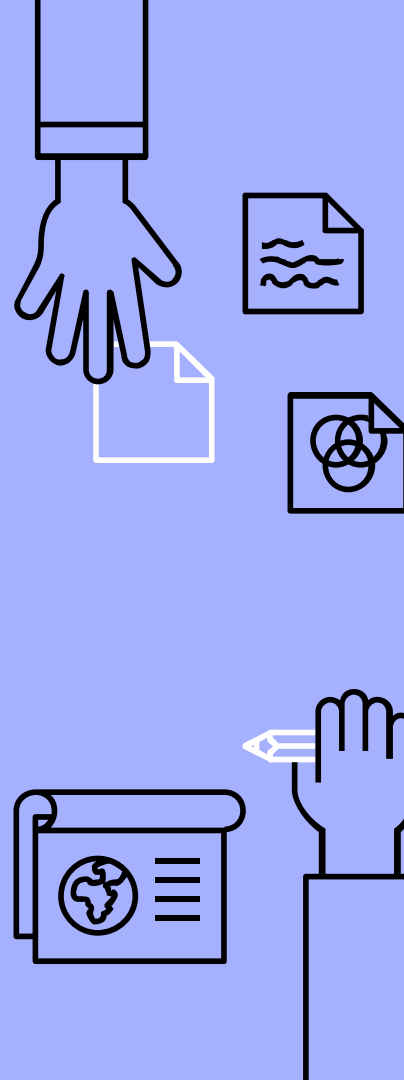
**Likelihood**  
How probable is the evidence  
given that our hypothesis is true?

**Prior**  
How probable was our hypothesis  
before observing the evidence?

$$P(H | e) = \frac{P(e | H) P(H)}{P(e)}$$

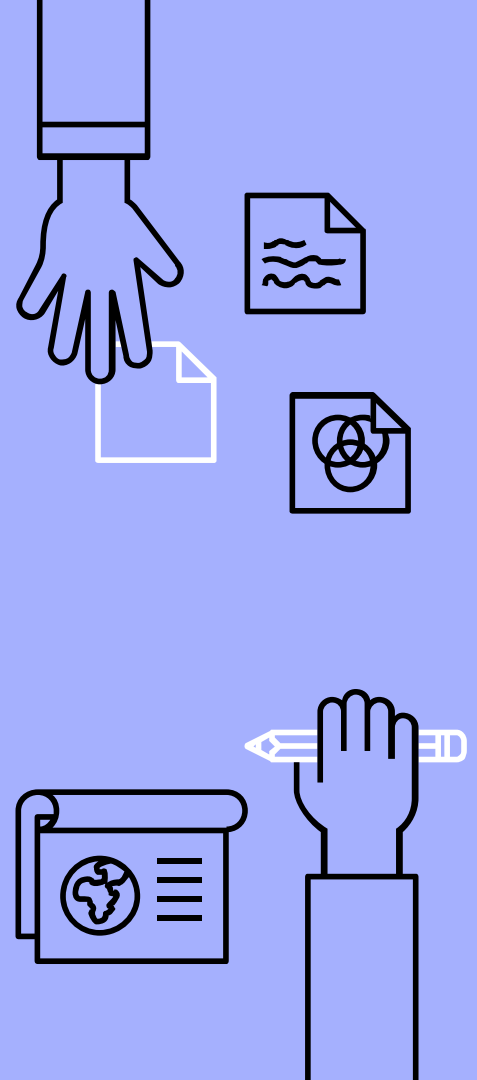
**Posterior**  
How probable is our hypothesis  
given the observed evidence?  
(Not directly computable)

**Marginal**  
How probable is the new evidence  
under all possible hypotheses?  
 $P(e) = \sum P(e | H_i) P(H_i)$



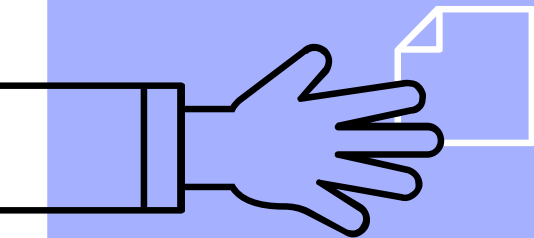
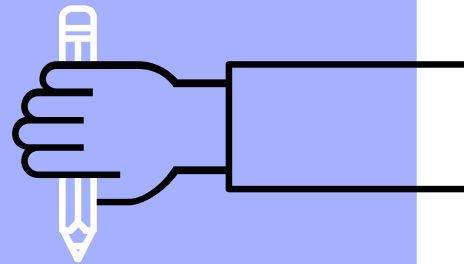
# Naive Bayes

- ▶ Reduced levels by putting the numerical values in categorical interval and everything else was in factors
- ▶ Installed the package e1071
- ▶ We ran the test with the categorical variables





KNN



# KNN

Accuracy: 0.9883464

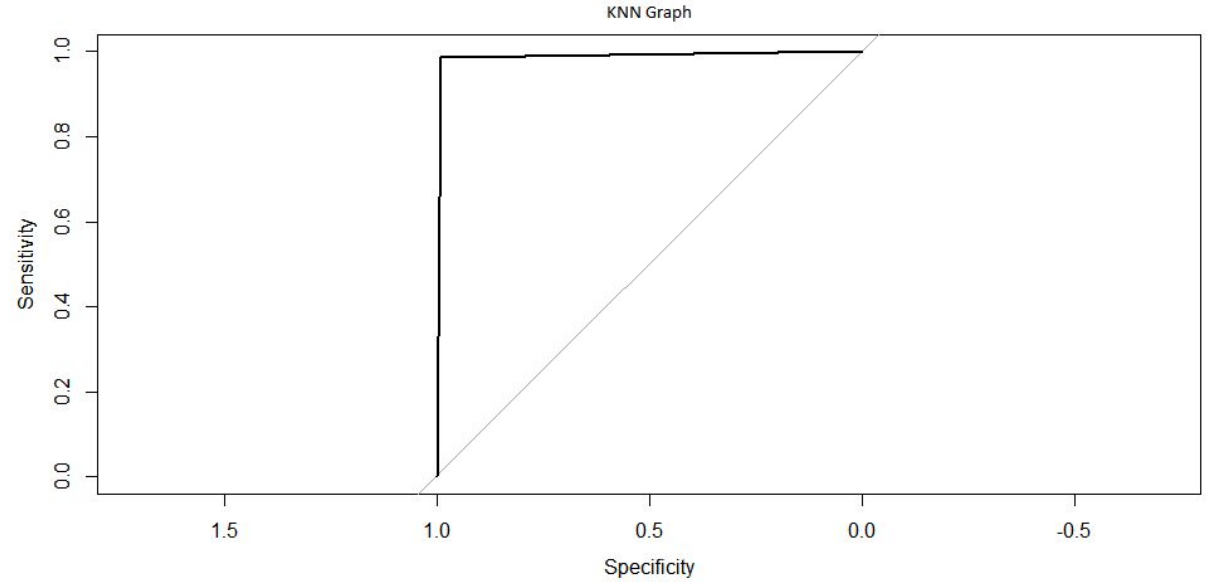
Sensitivity: 0.9870273

Specificity: 0.9869397

F1: 0.9869835

AUC: 0.9884

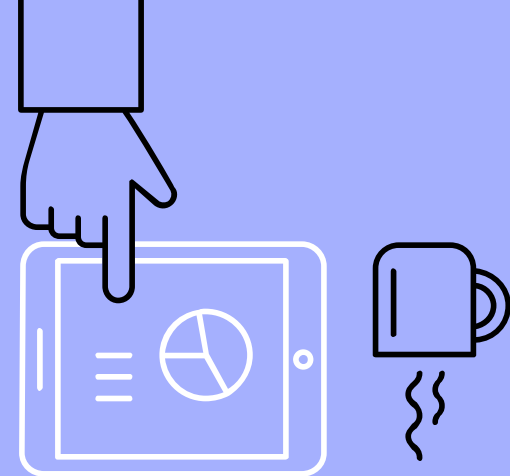
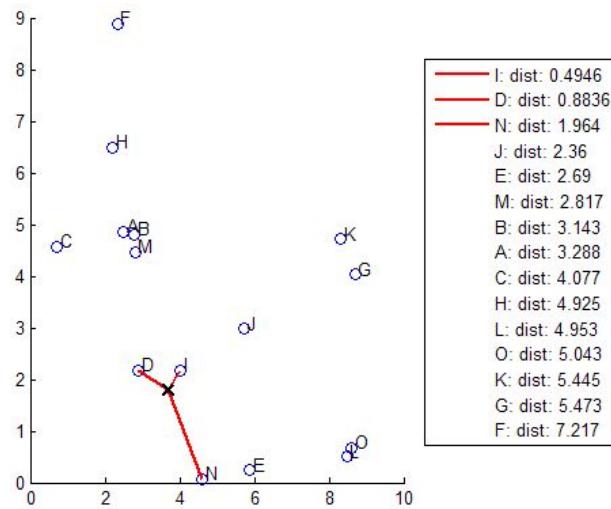
	False	True
0	8917	118
1	93	8978



# KNN

- ▶ Just like the other models we started by reducing dimensions
- ▶ Then we found out what would be the best k value
- ▶ Then we ran our model

Here is a picture we found  
To better explain the model



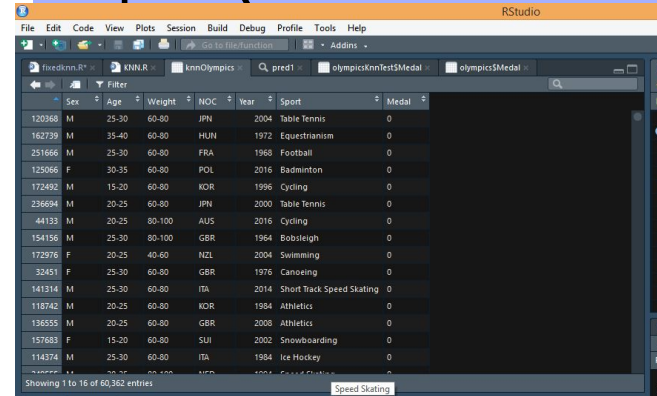


# Best Model??

**KNN** was the **best model** for our data by far. We used the variables: **Sex, Weight, Age, NOC, Year, and Sport**. We further broke down the Weight and Age variables by assigning them to intervals.

Next best was our **GLM** with the predictors: **Height, Event, Sex, Age, NOC**. The sensitivity and specificity were pretty average around 70%

Lastly, the Naive Bayes fell in last place because of it's low accuracy. We used the predictors **Sex, Age, Weight, Height, NOC Year, and Sport**.



The image shows a screenshot of the RStudio interface. The main window displays a data table with columns: Sex, Age, Weight, NOC, Year, Sport, and Medal. The table contains 16 rows of data, showing athletes from various countries (JPN, HUN, FRA, POL, KOR, AUS, GBR, NZL, ITA, SUI, JPN, ITA, etc.) and their participation in different sports (Table Tennis, Equestrianism, Football, Badminton, Cycling, Bobsleigh, Swimming, Short Track Speed Skating, Athletics, Snowboarding, Ice Hockey, etc.). The Medal column shows 0 for all entries. The status bar at the bottom indicates 'Showing 1 to 16 of 60,362 entries' and 'Speed Skating'.

	Sex	Age	Weight	NOC	Year	Sport	Medal
120368	M	25-30	60-80	JPN	2004	Table Tennis	0
162739	M	35-40	60-80	HUN	1972	Equestrianism	0
251666	M	25-30	60-80	FRA	1968	Football	0
725066	F	30-35	60-80	POL	2016	Badminton	0
172492	M	15-20	60-80	KOR	1996	Cycling	0
236694	M	20-25	60-80	JPN	2000	Table Tennis	0
44133	M	20-25	80-100	AUS	2016	Cycling	0
154156	M	25-30	80-100	GBR	1964	Bobsleigh	0
172976	F	20-25	40-60	NZL	2004	Swimming	0
32451	F	25-30	60-80	GBR	1976	Canoeing	0
141314	M	25-30	60-80	ITA	2014	Short Track Speed Skating	0
118742	M	20-25	60-80	KOR	1984	Athletics	0
136555	M	20-25	60-80	GBR	2008	Athletics	0
157683	F	15-20	60-80	SUI	2002	Snowboarding	0
114374	M	25-30	60-80	ITA	1984	Ice Hockey	0



# THANKS!

## Any questions?

You can find us at:

@smeragora on GitHub

And

@ShwinyG on GitHub

