

**THE ANTI-  
HEART  
DISEASE  
WARRIORS**

MASON CHEN BLACK BELT, STANFORD OHS

1<sup>ST</sup> PLACE BEST CONTRIBUTED PAPER, 2018 JMP DISCOVERY SUMMIT, CARY NC

# Project Scope and Presentation Flow

- Many people like chocolate, but have some concerns that chocolate is unhealthy.
- Some people who have **heart diseases** might need to eat chocolate, but do not know which one to eat.

1. Anti-Oxidant Science  
Literature Research

4. Missing Value Neural  
Imputation of Cocoa%

3. Clustering Chocolate  
Types

2. Clustering Nutritions  
& Science

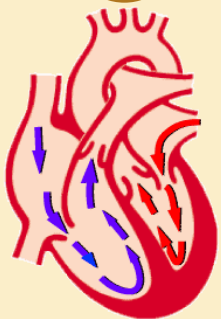
5. DSD Optimization of  
Neural Setting



Ingredients and Nutritions

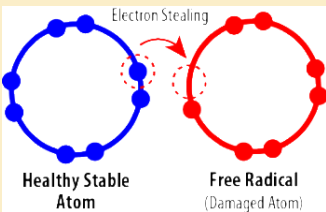
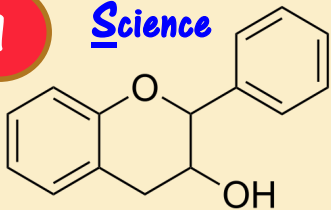


2 Clustering Nutritions  
(Artificial Intelligence)

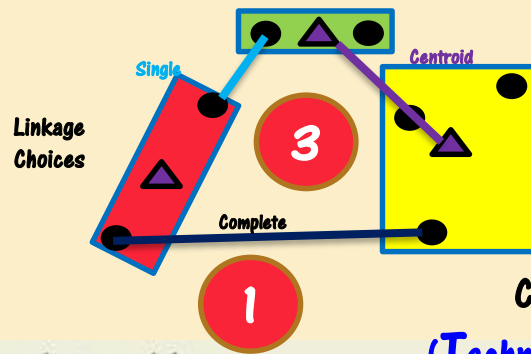


Anti-Oxidant

1 Science



# ST EAMS DIAGRAM



3 Clustering Products  
(Artificial Intelligence)  
Clustering Algorithm  
(Math)

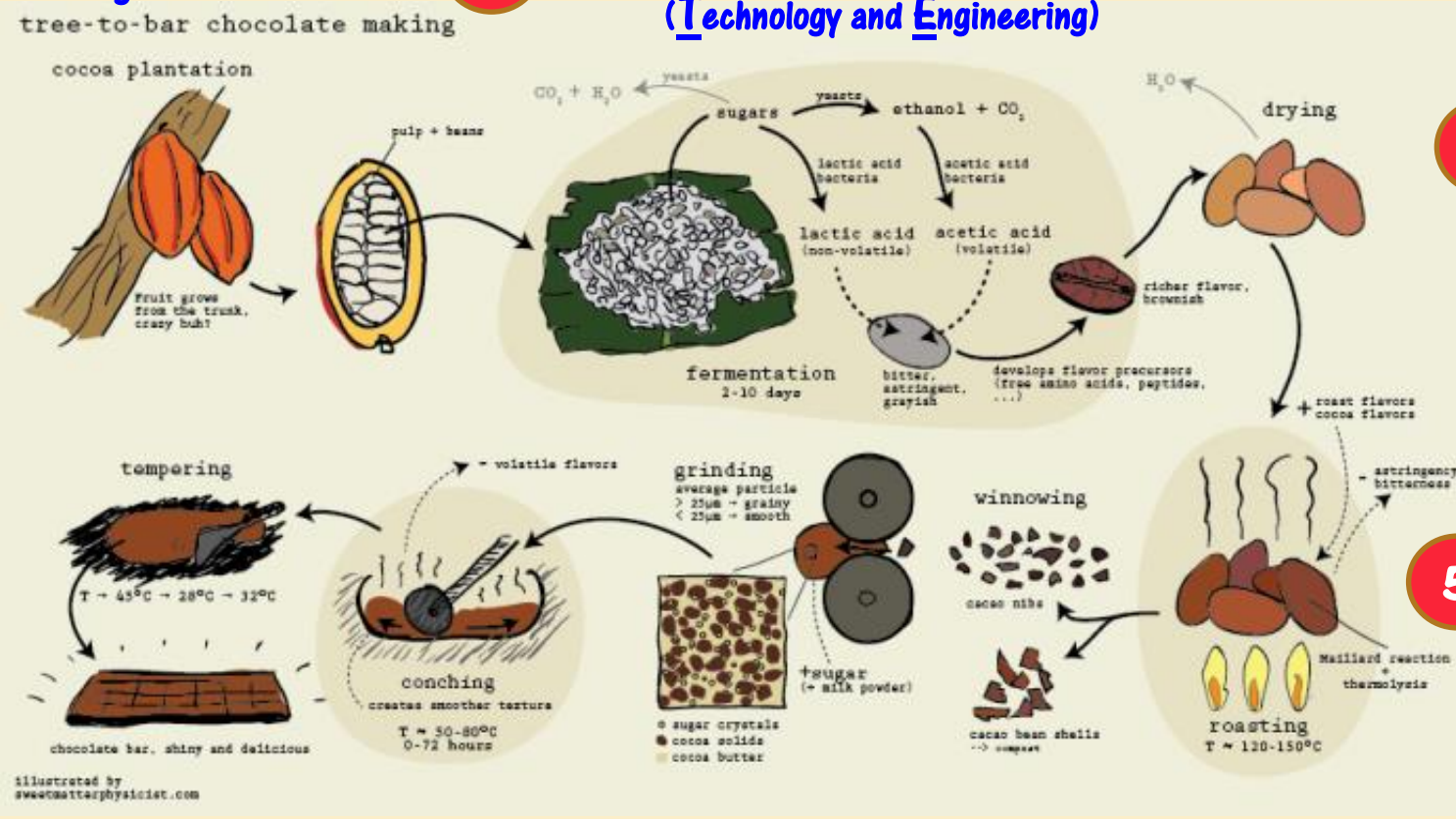


Chocolate Products

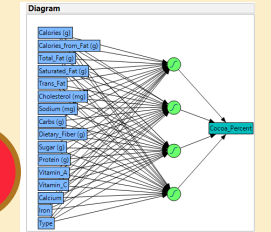


Chocolate Process

(Technology and Engineering)



4



Neural Imputation  
(Artificial Intelligence)

Power Analysis

Significance Level 0.05

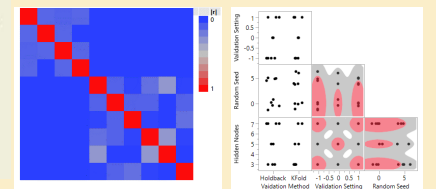
Anticipated RMSE 1

Term	Anticipated Coefficient	Power
Intercept	1	0.975
Validation Method	1	0.97
Validation Setting	1	0.925
Random Seed	1	0.97
Hidden Nodes	1	0.925

Apply Changes to Anticipated Coefficients

5

DSD Optimization  
(Statistics)



# 2015-2018 "STEAMS" Journey

S T E M → T E A M S → S T E A M S

Develop Math & Science  
Foundation  
(Stanford OHS)

Math, Physics, Biology, Chemistry,  
JAVA, **Statistics** Literature  
Research/Writing

Fun, Real

Certify 6 Professional  
Certificates:

IBM SPSS Statistics  
IBM Modeler DA/DM (IBM000129876)  
IASSC YB/GB/BB (GR764000541MC)  
JMP **Statistical Thinking** (2018 Goal)  
JMP **DOE** (2019 Goal)  
JMP **Script Specialist** (2019 Goal)

Hands-On

Enhance STEAMS Skills  
**JMP/Pro**, Python  
Latex Paper Proceedings  
Oral/Poster Presentations  
Team Building

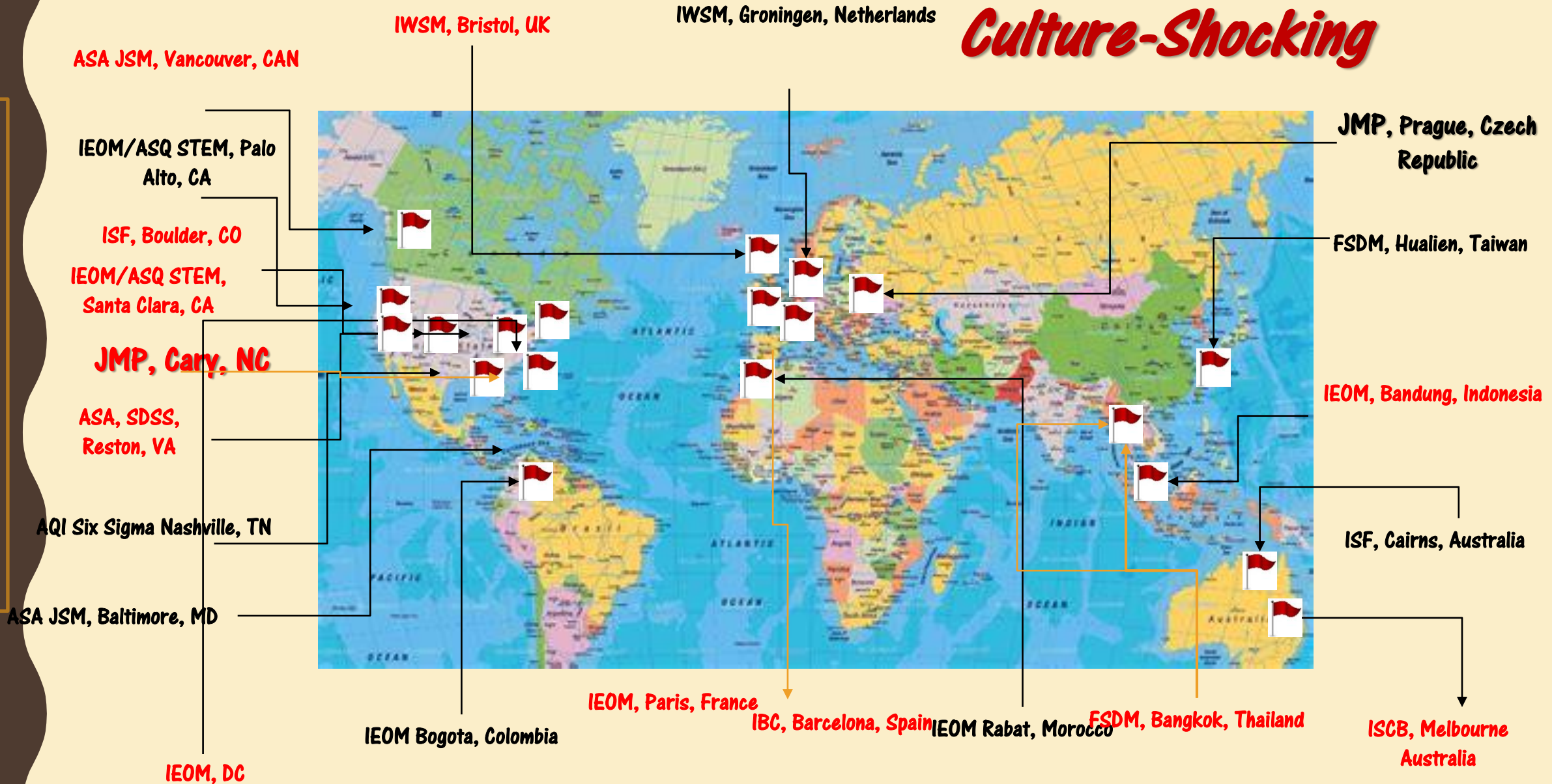
Learning "STEAMS" techniques help motivate school learning on project-based and practical way

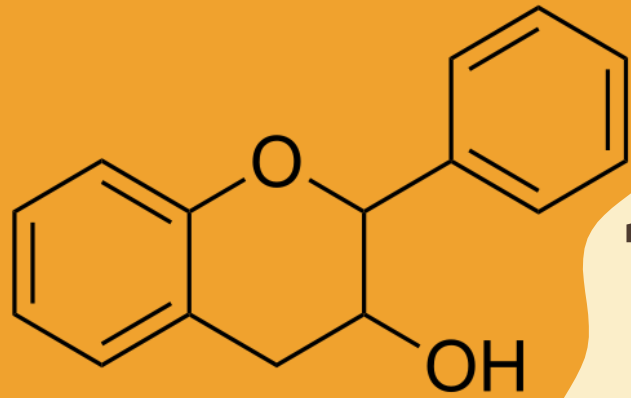


# Global Vision Leadership: 2017-2018 Conferences

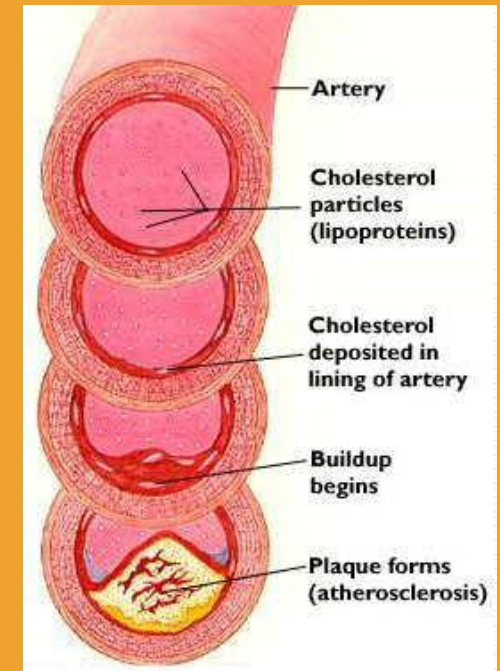
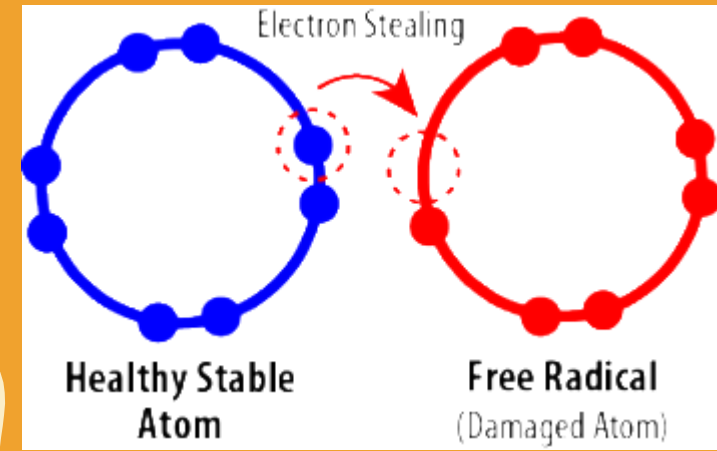
## Culture-Shocking

STEAMS



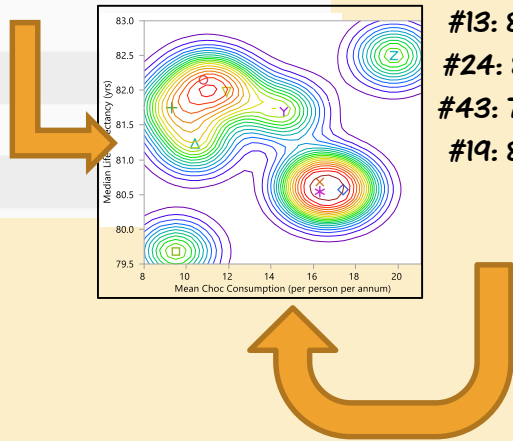
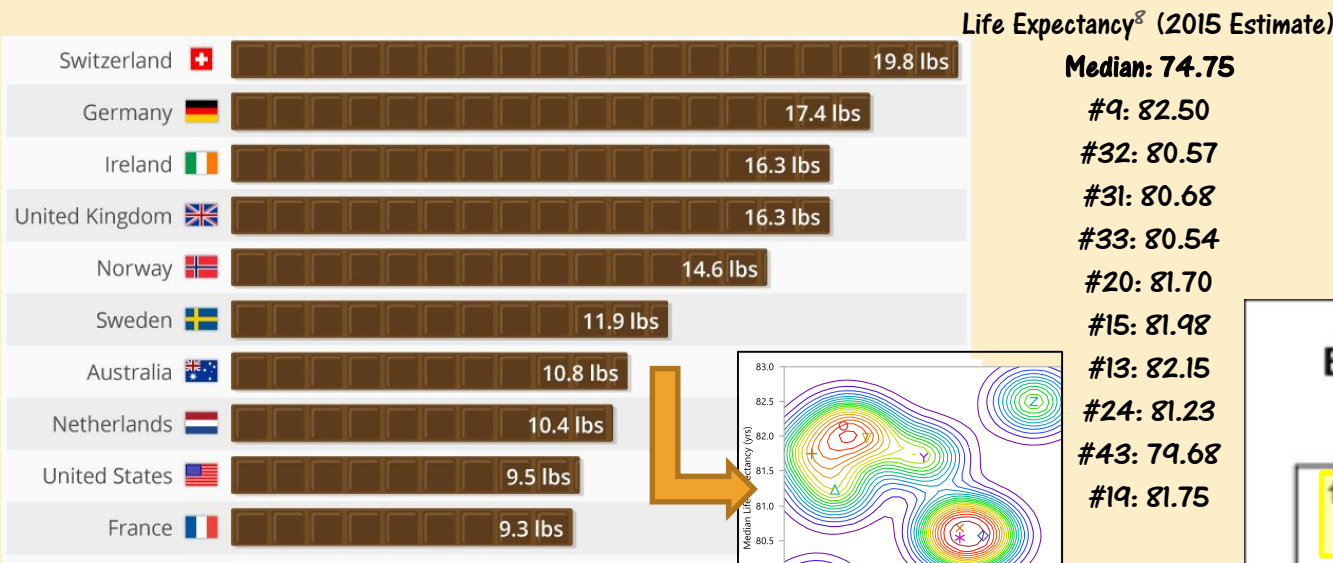


# 1. Chocolate Anti-Oxidant Science



# IS EATING CHOCOLATE UNHEALTHY?

Chocolate has not been proven harmful.



JMPI3 >> Analyze >> Fit Y by X >> Nonpar Density

Anti-Oxidant Capacity/gram

**Estimates of Antioxidant Capacity for Selected Foods**  
(micromole TE per household measure and grams)

	0	2000	4000	6000
1 sm, 149 g Apple, Red Delicious, w/skin				6370
1 oz, 28 g Chocolate, Dark				5903
1/2 c, 87 g Plums, dried				5700
5 fl oz, 147 g Wine, red				5693
1/2 med, 60 g Artichokes, Ocean Mist, boiled				5650
1 oz, 28 g Pecans				5023
1/2 c, 74 g Blueberries, fresh			4848	
1 oz, 28 g Walnuts, English		3791		
1/2 c, 83 g Strawberries, sliced	2969			
1 med, 114 g Sweet potato, baked	2411			

Source: Calculated from *Oxygen Radical Absorbance Capacity of Selected Foods, 2007*  
USDA-Agricultural Research Service  
(www.ars.usda.gov/nutrientdata/ORAC)

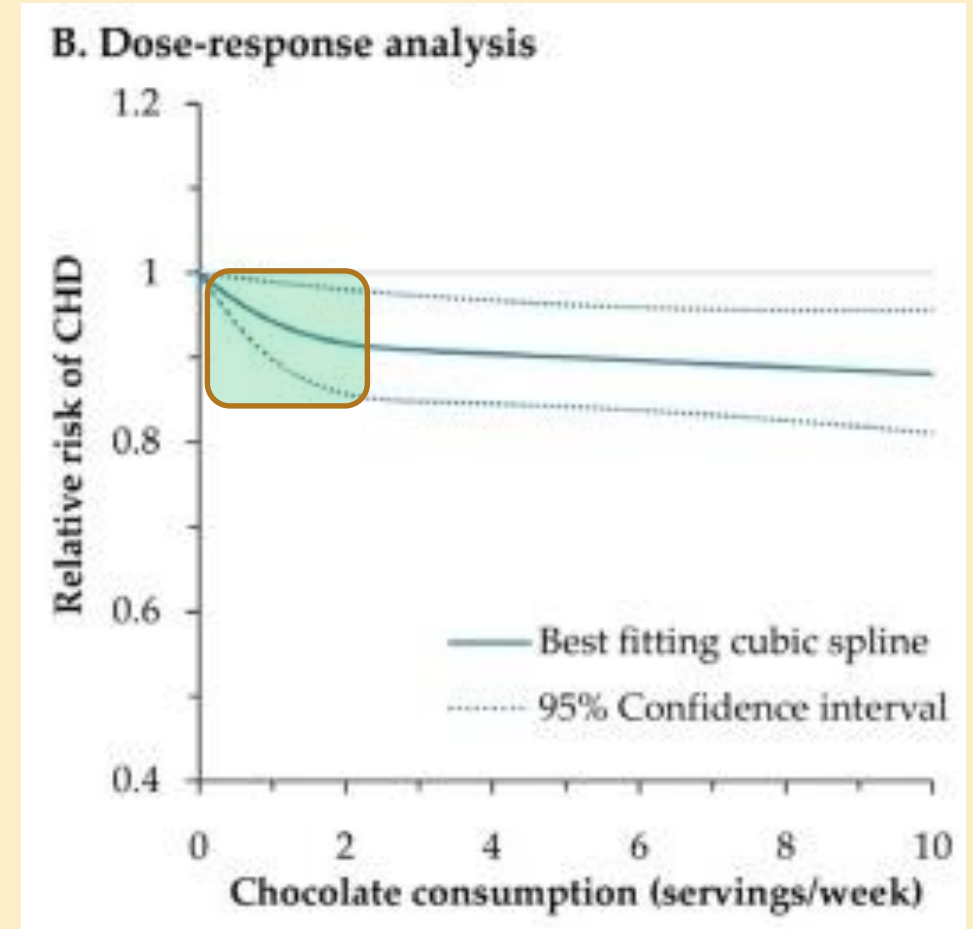
Dark chocolate is a powerful source of antioxidant. If chocolate's serving size is equal to that of an apple, it has the highest amount of antioxidant.



# CHOCOLATE & ATRIAL FIBRILLATION (AF)

Lower Cardiovascular Heart Disease (CHD) risk if taking 2 Chocolate servings per week (1 serving = 30 g)

- Chocolate may be inversely associated with AF
- Dark chocolate may be a healthy snacking option
- AF = Atrial Fibrillation (a cardiovascular disease)
- Next, how Chocolate can reduce CHD risk and AF associated cardiovascular disease



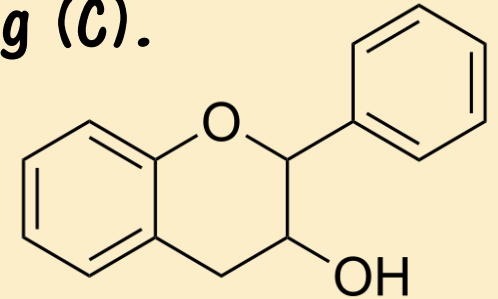
<https://heart.bmj.com/content/103/15/1163>

<https://www.bmj.com/content/343/bmj.d4488>

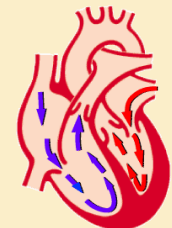


# FLAVONOIDS SCIENCE & STRUCTURE

- Flavonoids are the most abundant polyphenols in human diet that have antioxidant properties.
- Flavonoids have the general structure of a 15- carbon skeleton C6-C3-C6.
  - Consists of two phenyl rings (A and B) and a heterocyclic ring (C).

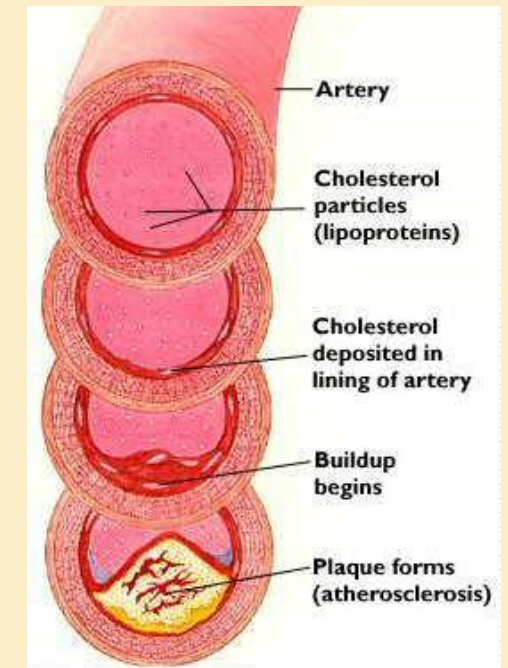
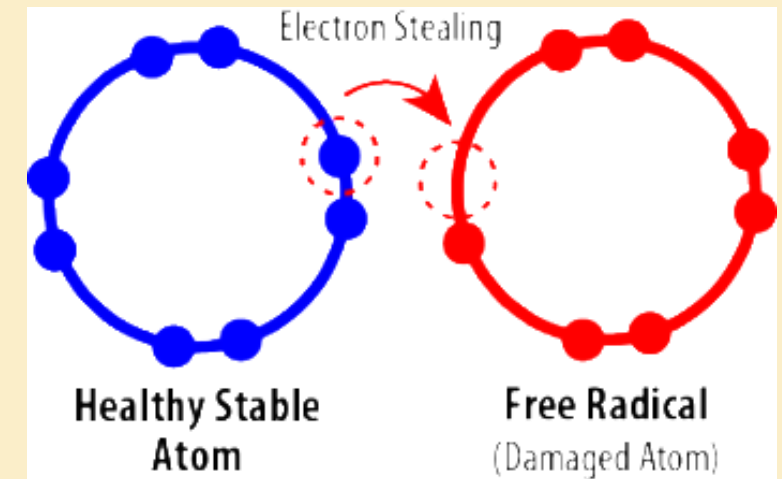


- There are seven different types of flavonoids based on its chemical structure:
  - Flavones, flavanol, flavanones, isoflavones, anthocyanidins, chalcones, catechins
- Chocolate flavonoids are flavanols which can promote healthy blood flow from head to toe.



# FREE RADICALS AND ANTIOXIDANTS

- Free radicals are atoms with odd number of electrons
  - **Antioxidants reduce free radical formation**
  - **Reactive free radicals causes cells mal-function**
  - **Excess free radicals damages blood vessel**
- **After the oxidation of free radicals, LDL (Low-density Lipoprotein) can cause CVD (Cardiovascular Disease)**
  - **The oxidized components attract macrophages which absorb & deposit Cholesterol**



# DARK CHOCOLATE LITERATURE RESEARCH

- **Benefits:**

- A lot of soluble fiber
- A lot of minerals: iron, magnesium, copper, manganese, potassium, phosphorus, zinc, selenium
- Powerful source of antioxidant
- Improve blood flow and lower blood pressure
- Increases HDL (good cholesterol) and decreases LDL (bad cholesterol)
- Lower risk of cardiovascular disease (CVD)
- Improve brain function<sup>1</sup>

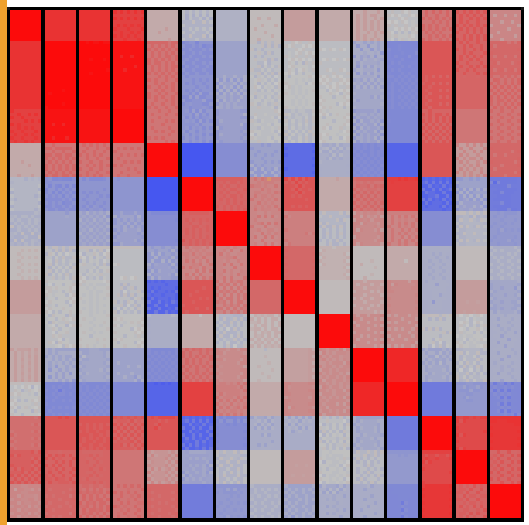
- **Concerns:**

- Causes migraines
- Increases chance of kidney stones
- Side effects from caffeine such as irregular heartbeat





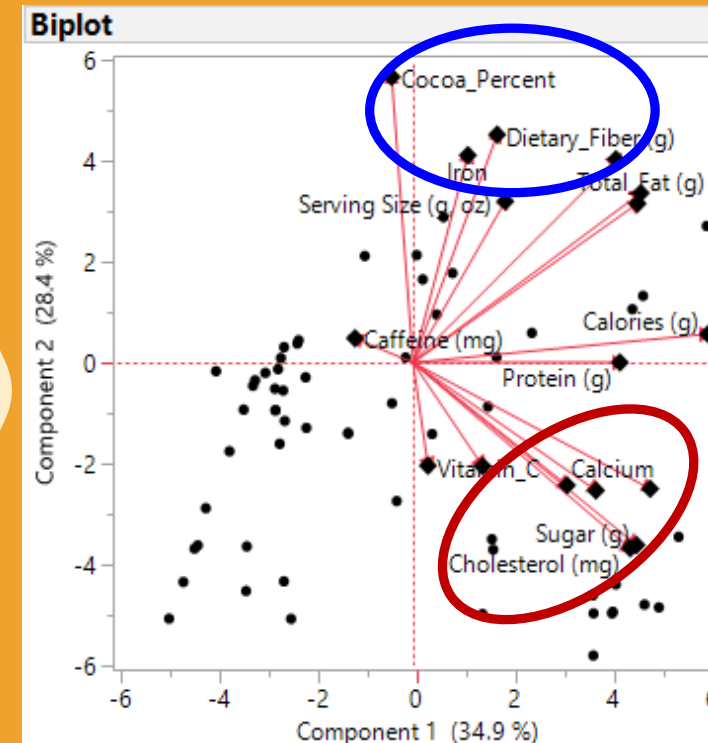
(1) JMP 13 >> Analyze >>  
Distribution



(2) JMP 13 >> Analyze >>  
Multivariate Methods >>  
Multivariate

## 2. Clustering Nutritions & Science

(3) JMP 13 >> Analyze >> Clustering >>  
Cluster Variables



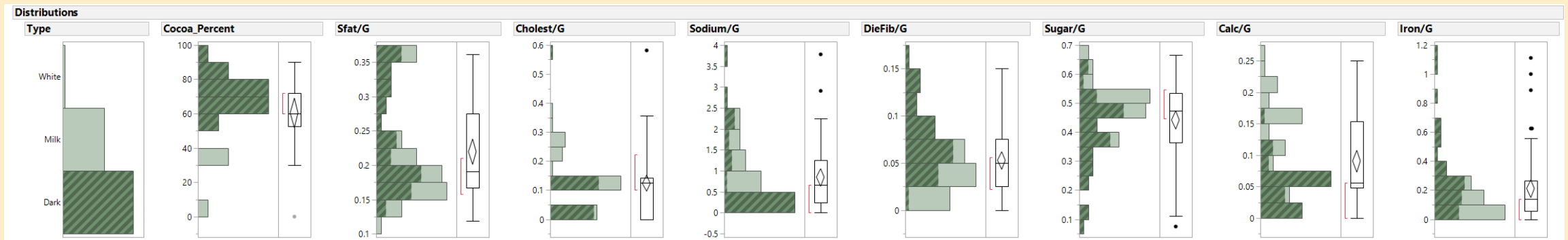
(4) JMP 13 >> Analyze >> Multivariate  
Methods >> Principle Components >> Bi-Plot

• Mason C., (2018 July), "Choose Healthy Chocolate", IEOM Europe Paris Proceedings, 434-441

# (1) CHOCOLATE NUTRITION DISTRIBUTION

JMP 13 >> Analyze >>  
Distribution

- 60+ Chocolate nutrition data collected from “Target” store.
- The quantities of the eight most critical ingredients were analyzed.



Most Dark Chocolate (**Qualitative** Clustering Criteria) has:

- 1<sup>st</sup> Cluster: higher Cocoa percent, Dietary Fiber and Iron
- 2<sup>nd</sup> Cluster: lower Cholesterol, Calcium, and Sugar

Chocolate Product Nutrition data has indicated that Dark Chocolate is healthier than the Milk and White Chocolate

# (2) DARK CHOCOLATE CORRELATION

- **1<sup>st</sup> Cluster:** Sugar and Cocoa\_Percent have a negative correlation of **-0.9162**.
- **2<sup>nd</sup> Cluster:** Dietary Fiber and Iron have a positive correlation of **0.7722**.

A 100 gram bar of dark chocolate with 70-85% cocoa contains (1):

- 11 grams of fiber.
- 67% of the RDA for Iron.
- 58% of the RDA for Magnesium.
- 89% of the RDA for Copper.
- 98% of the RDA for Manganese.
- It also has plenty of potassium, phosphorus, zinc and selenium.

Any other better way to cluster nutritions?

Correlations	Pair-Wise Pearson Correlation							
	Cocoa_Percent	Sfat/G_1	Cholest/G_1	Sodium/G_1	DieFib/G_1	Sugar/G_1	Calc/G_1	Iron/G_1
Cocoa_Percent	1.0000	0.5291	-0.3114	-0.0583	0.5482	-0.9162	0.2625	0.4597
Sfat/G_1	0.5291	1.0000	-0.1980	0.0184	0.0341	-0.7068	0.4161	0.0687
Cholest/G_1	-0.3114	-0.1980	1.0000	0.0302	-0.3666	0.3333	0.1732	-0.3304
Sodium/G_1	-0.0583	0.0184	0.0302	1.0000	-0.1344	0.0462	0.1667	-0.1862
DieFib/G_1	0.5482	0.0341	-0.3666	-0.1344	1.0000	-0.5804	-0.0207	0.7722
Sugar/G_1	-0.9162	-0.7068	0.3333	0.0462	-0.5804	1.0000	-0.3696	-0.4669
Calc/G_1	0.2625	0.4161	0.1732	0.1667	-0.0207	-0.3696	1.0000	-0.1037
Iron/G_1	0.4597	0.0687	-0.3304	-0.1862	0.7722	-0.4669	-0.1037	1.0000

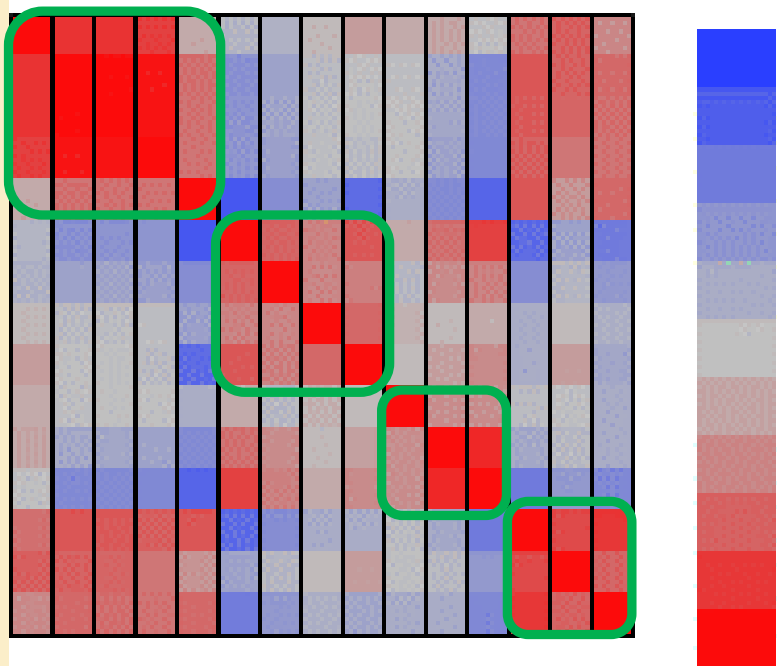
JMP 13 >> Analyze  
>> Multivariate  
Methods >>  
Multivariate



# (3) VARIABLE CLUSTERING

JMP 13 >> Analyze >> Clustering >> Cluster Variables

Color Map on Correlations



Cluster Members		Signal	Noise	S-N Ratio
Cluster	Members	RSquare with Own Cluster	RSquare with Next Closest	1-RSquare Ratio
1	Calories (g)	0.789	0.314	0.308
1	Calories_from_Fat (g)	0.976	0.456	0.044
1	Total_Fat (g)	0.977	0.426	0.04
1	Saturated_Fat (g)	0.935	0.361	0.101
2	Cocoa_Percent	0.742	0.366	0.406
2	Cholesterol (mg)	0.811	0.387	0.309
2	Vitamin_A	0.505	0.126	0.566
2	Vitamin_C	0.412	0.016	0.598
2	Calcium	0.726	0.079	0.297
3	Sodium (mg)	0.345	0.013	0.664
3	Carbs (g)	0.876	0.185	0.152
3	Sugar (g)	0.874	0.416	0.216
4	Dietary_Fiber (g)	0.888	0.403	0.187
4	Protein (g)	0.73	0.358	0.421
4	Iron	0.803	0.269	0.269

Clustering Nutritions can interpret the relevant Chocolate Science insight well:

Cluster 1: the higher the saturated fat, the higher the total fat, and the higher the calories.

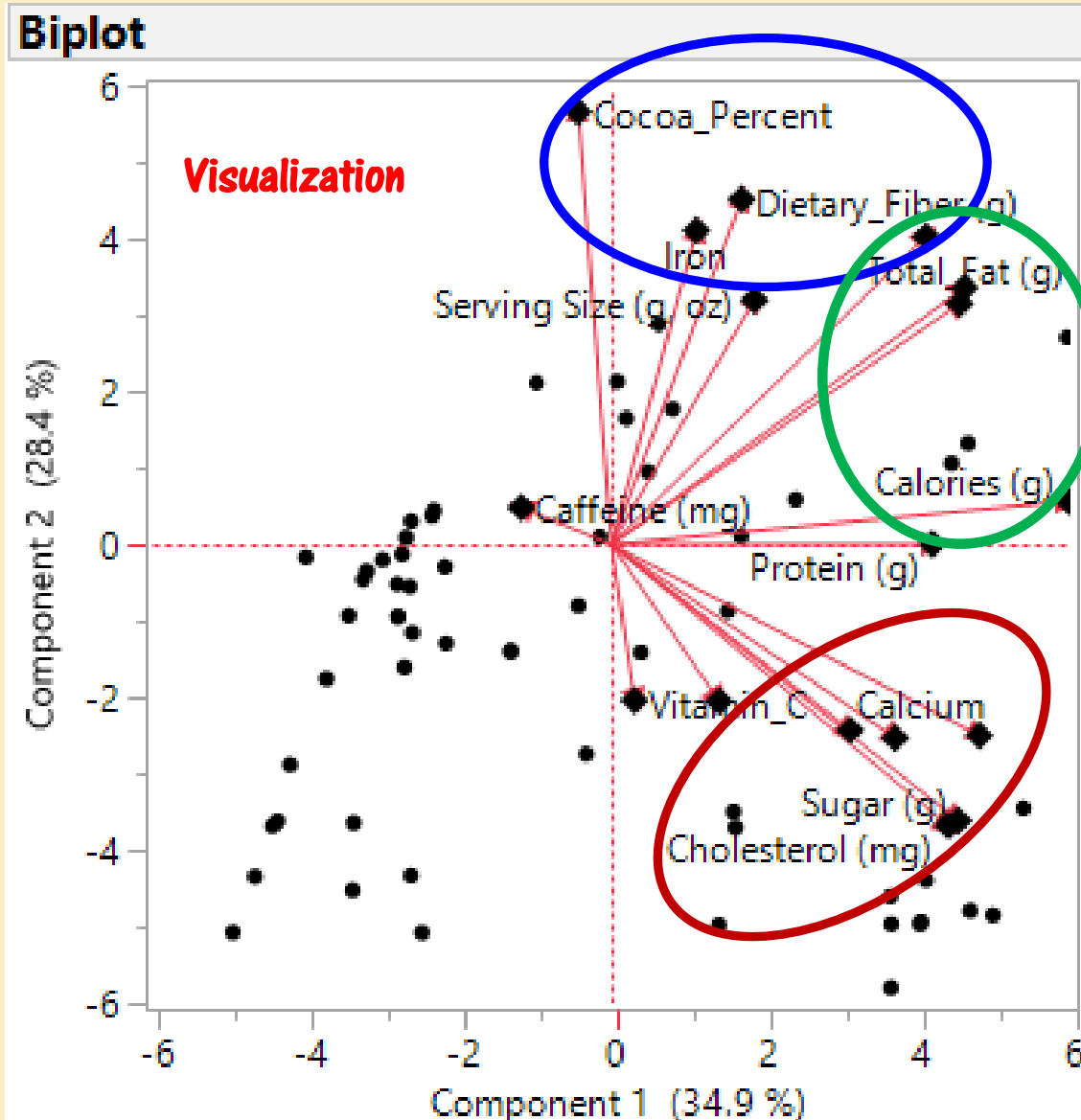
Cluster 2: Calcium/Cholesterol, and Cocoa percent have a negative correlation.

Cluster 3: the higher the sugar, the higher the carbohydrates.

Cluster 4: Iron and dietary fiber are positively correlated.

Common Sense

# (4) Principle Component Bi-Plot



**1<sup>st</sup> Cluster:** Cocoa Percent, Dietary Fiber, and Iron are near each other (Higher for Dark Chocolate)

**2<sup>nd</sup> Cluster:** Total Fat, Saturated Fat, and Calories

**3<sup>rd</sup> Cluster:** Calcium, Sugar, and Cholesterol are near each other (Higher for Milk/White Chocolate)

JMP 13 >> Analyze >> Multivariate Methods >> Principle Components >> Bi-Plot

# Comparing Four Clustering Methods

Platform	Criteria	1st Cluster	2nd Cluster	3rd Cluster	4th Cluster
Interactive Distribution	Qualitative	Cocoa%, Dietary Fiber, Iron	Cholestrol, Sugar		
Multivariate Correlation	Quantitative	Dietary Fiber and Iron	Cocoa% and Sugar		
Clustering Variables	Quantitative	Saturated Fat, Total Fat, Calories	Cholesterol, Calcium, Cocoa%	Sugar, Carbohydrates	Iron, Dietary Fiber
Principle Component Bi-Plot	Quantitative	Cocoa %, Dietary Fiber, Iron	Saturated Fat, Total Fat, Calories	Calcium, Sugar, Cholesterol	

- Four different clustering methods show similar clustering patterns
- Clustering “**Statistics** and **Engineering**” results match Chocolate “**Science** and **Technology**” Literature Research well (**STEAMS**).



JMP 13 >> Analyze >> Clustering >>  
Hierarchical Cluster  
JMP 13 >> Analyze >> Distribution

JMP 13 >> Analyze >> Clustering >>  
Hierarchical Cluster >> Clustering Distance  
Method

# 3. Clustering Chocolate Types

JMP 13 >> Analyze >>  
Distribution

JMP 13 >> Analyze >>  
Clustering >> Hierarchical  
Cluster >> Constellation Plot

JMP 13 >> Analyze >> Clustering >>  
Hierarchical Cluster >> Column  
Summary

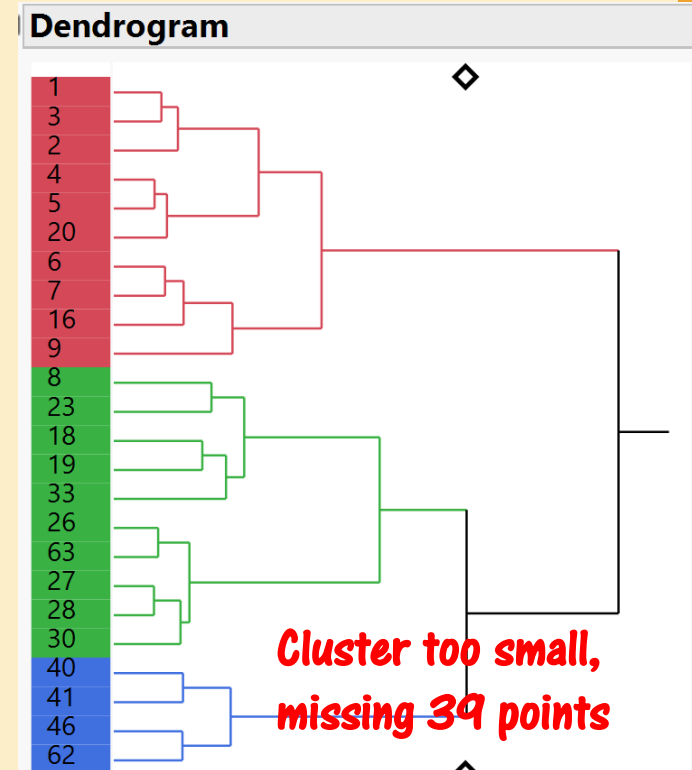
JMP 13 >> Analyze >> Multivariate  
Methods >> Principle Components  
>> Eigenvalues

- Mason C., (2018 Dec.), "Statistics Application on the Study of Chocolate Science with Heart Disease", ASA SDSS Proceeding

# CLUSTERING PRODUCTS

Objective: find a way to identify healthy chocolate products for **Heart Disease** patients.

- Use hierarchical clustering to cluster chocolate products
- All Milk and white chocolate form the third cluster while dark chocolate split between the first and second cluster.
- Why are there two clusters for dark chocolate (why not one cluster for each Chocolate Type)?

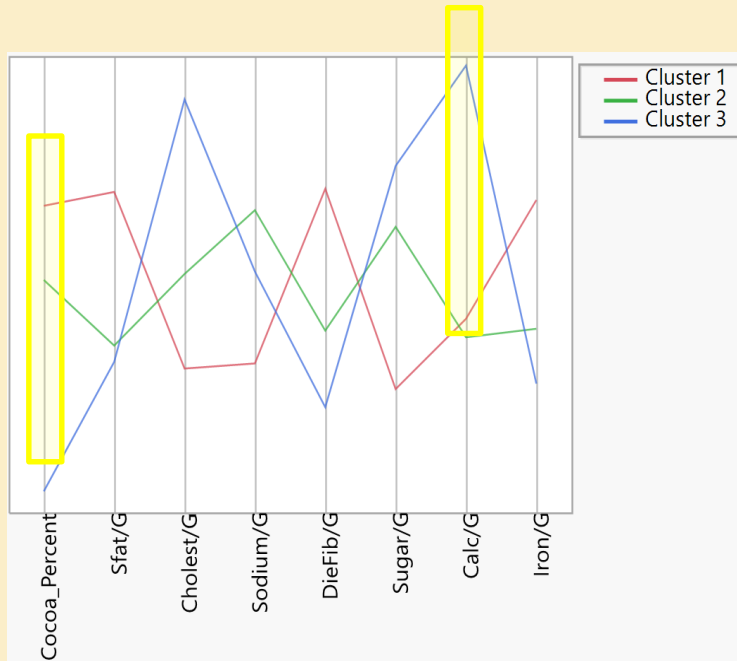


JMP 13 >> Analyze >> Clustering >> Hierarchical Cluster

JMP 13 >> Analyze >> Distribution

# PRINCIPLE CLUSTERING DECIDING FACTORS

JMP 13 >> Analyze >> Clustering >> Hierarchical Cluster >> Column Summary



## Column Summary

Column	RSquare	.2	.4	.6	.8
Cocoa_Percent	0.7788	[Green bar]			
Sfat/G	0.4908	[Grey bar]			
Cholest/G	0.6978	[Grey bar]			
Sodium/G	0.3949	[Grey bar]			
DieFib/G	0.5788	[Grey bar]			
Sugar/G	0.6642	[Grey bar]			
Calc/G	0.7727	[Green bar]			
Iron/G	0.4351	[Grey bar]			

- 1<sup>st</sup> Cluster: Dark Chocolate, High Cocoa%, and Low Calcium, **Most Healthy?**
- 2<sup>nd</sup> Cluster: Dark Chocolate, Medium Cocoa%, and Low Calcium.
- 3<sup>rd</sup> Cluster: Milk/White Chocolate, Low Cocoa%, and High Calcium

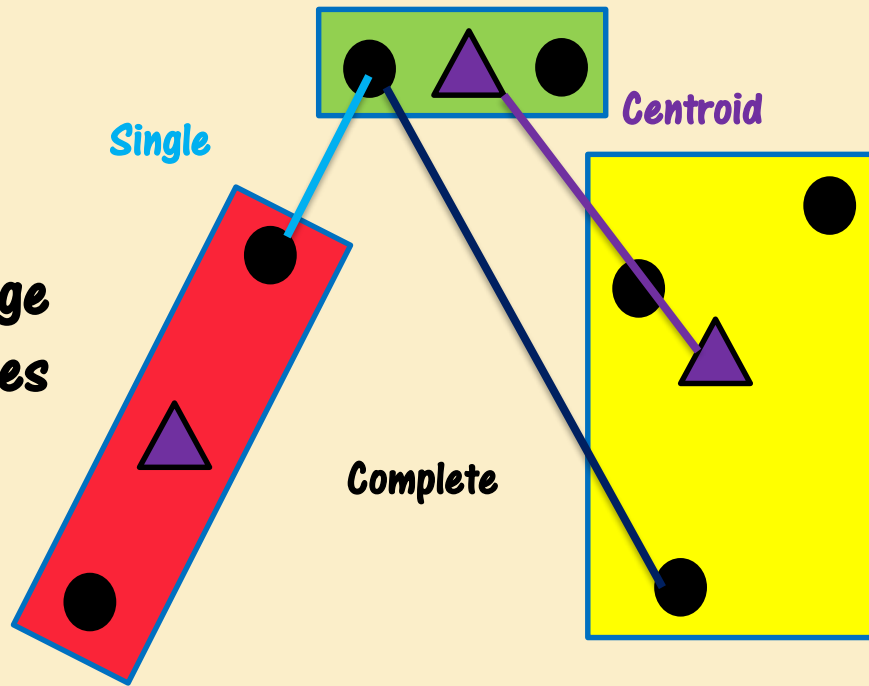


# CLUSTERING DISTANCE METHODS

JMP 13 >> Analyze >> Clustering >> Hierarchical Cluster

MATH

Linkage Choices



Clustering patterns dependent on the cluster number observations, cluster variance, and outlier

**Average Linkage** Distance for the average linkage cluster method is:

$$D_{KL} = \frac{\sum_{i \in C_K} \sum_{j \in C_L} d(x_i, x_j)}{N_K N_L} \quad \leftarrow \text{Average}$$

**Centroid Method** Distance for the centroid method of clustering is:

$$D_{KL} = \|\bar{x}_K - \bar{x}_L\|^2$$

**Ward's** Distance for Ward's method is: **Center-Center**

$$D_{KL} = \frac{\|\bar{x}_K - \bar{x}_L\|^2}{\frac{1}{N_K} + \frac{1}{N_L}} \quad \leftarrow \text{ANOVA (MS)}$$

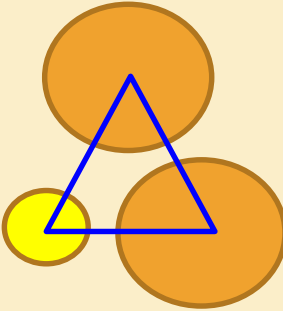
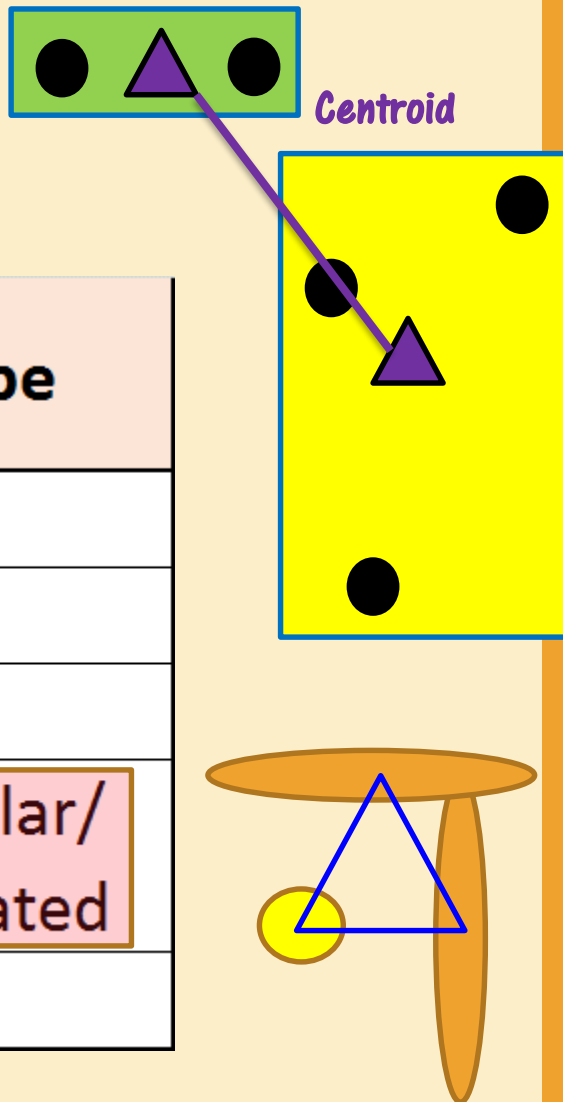
**Single Linkage** Distance for the single linkage cluster method is:

$$D_{KL} = \min_{i \in C_K} \min_{j \in C_L} d(x_i, x_j) \quad \leftarrow \text{Min}$$

**Complete Linkage** Distance for the Complete linkage cluster method is:

$$D_{KL} = \max_{i \in C_K} \max_{j \in C_L} d(x_i, x_j) \quad \text{Maximum}$$

# Selecting DISTANCE METHODS



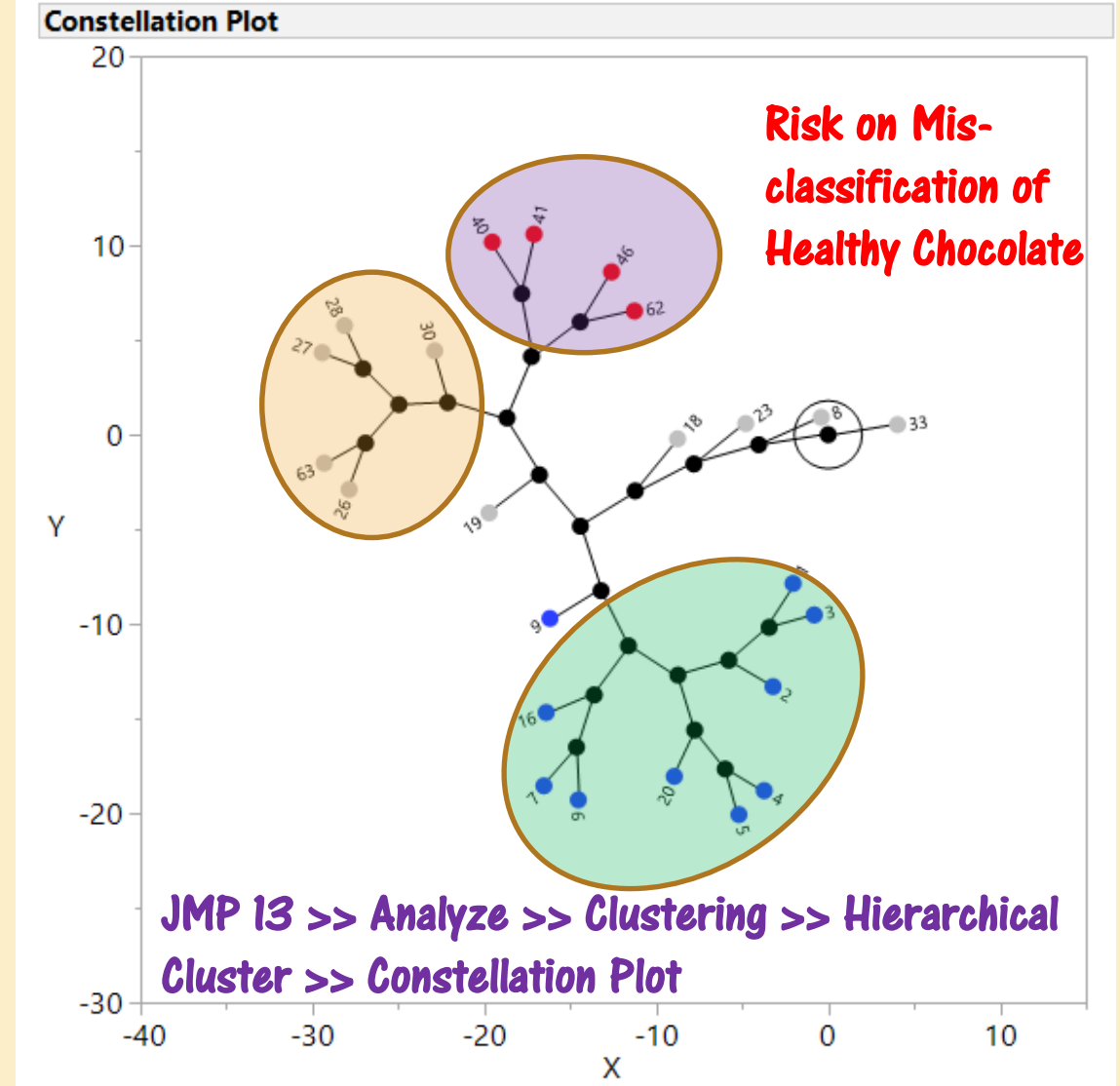
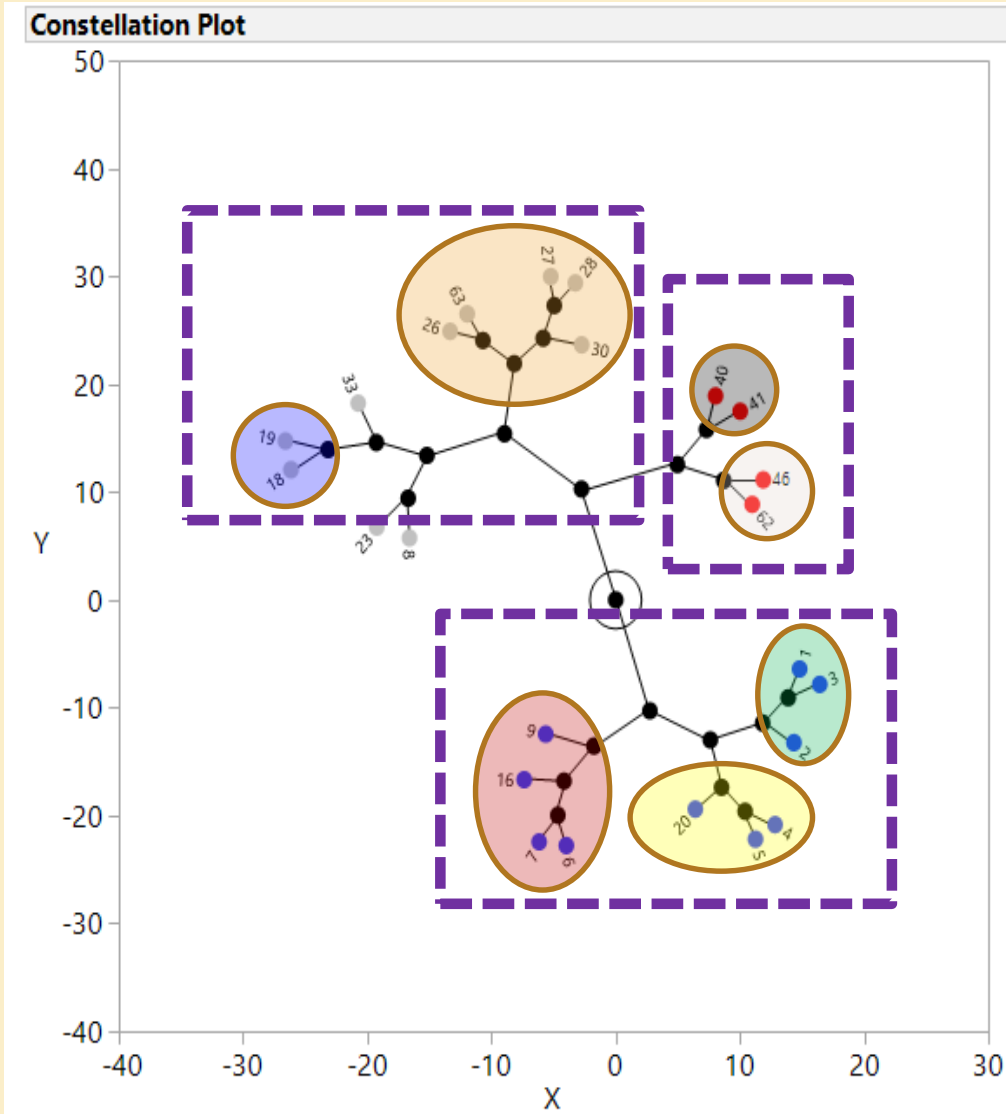
	Size/ Variance	Outliers	Shape
<b>Average</b>	Smaller		
<b>Centroid</b>		Robust	
<b>Ward</b>	Smaller	Sensitive	
<b>Single</b>	Larger		Irregular/ Elongated
<b>Complete</b>	Smaller	Moderate	

Depending on the data distribution, selecting an appropriate Clustering Distance algorithm is critical to Clustering Pattern Analysis

# WARD VS SINGLE METHOD (10 Clusters)

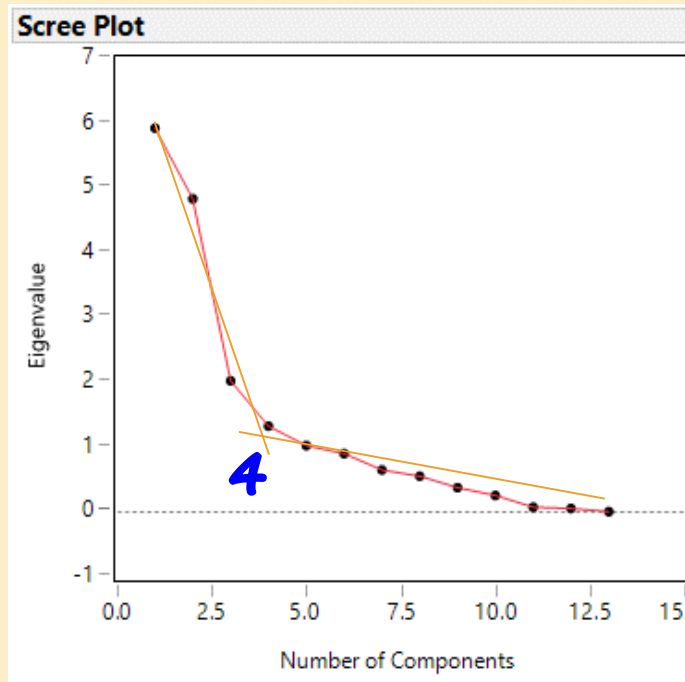
Ward (Join Smaller Observations)

Single (Join Larger Variances)



# DETERMINE NUMBER OF CLUSTERS

Clustering pattern result is highly dependent on the number of clusters



## Eigenvalues

Number	Eigenvalue	Percent	20	40	60	80	Cum Percent
1	5.9254	34.855					34.855
2	4.8336	28.433					63.288
3	2.0252	11.913					75.201
4	1.3256	7.797					82.998
5	1.0255	6.032					89.030
6	0.9010	5.300					94.330
7	0.6484	3.814					98.144

From both the scree plot and PCA eigenvalues (80% Pareto), we can pick **4** clusters

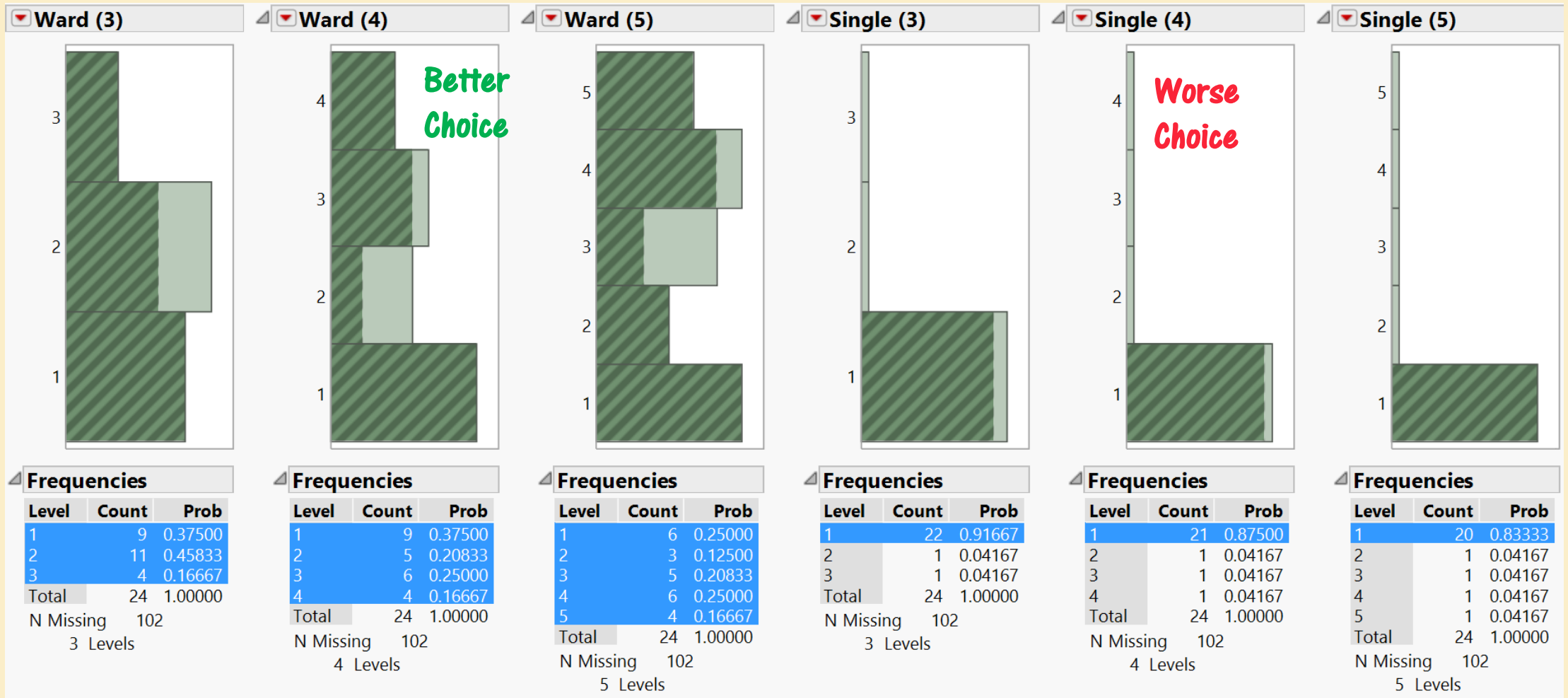
JMP 13 >> Analyze >> Multivariate Methods >> Principle Components >> Eigenvalues



# WARD VS SINGLE (3-5 CLUSTERS)

JMP 13 >> Analyze >> Distribution

STATISTICS



- Single does not show any significant difference between 3, 4, or 5 clusters
- Ward clusters become more similar in size with the higher the number of clusters

JMP 13 >> Analyze >> Clustering >>  
Hierarchical Cluster >> Missing Value  
Imputation

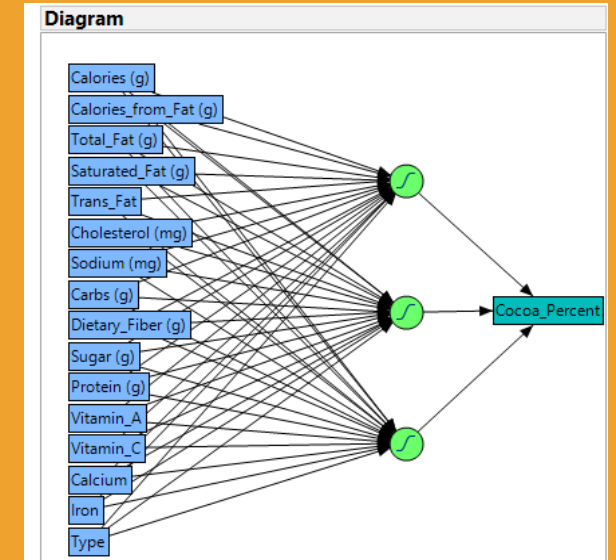
JMP 13 >> Analyze >>  
Screening >> Explore  
Missing Values

JMP 13 >> Analyze >>  
Predictive Modeling >>  
Partition

# 4. Missing Value Neural Imputation

- Mason C., "Neural Network Algorithm of Missing Value Imputation for Chocolate Science Research" submitted to SIAM SDMI9

JMP 13 >> Analyze >>  
Predictive Modeling >> Neural



JMP 13 >> Analyze >>  
Predictive Modeling >> Neural  
>> Diagram

# Explore MISSING VALUES

JMP 13 >> Analyze >> Screening >> Explore Missing Values

- Objective: among 63 commercial chocolate products, 39 have missed the Cocoa % information (most are Milk Chocolates)

**Missing Columns**

Show only columns with missing

Close

Column	Number Missing
Cocoa_Percent	39

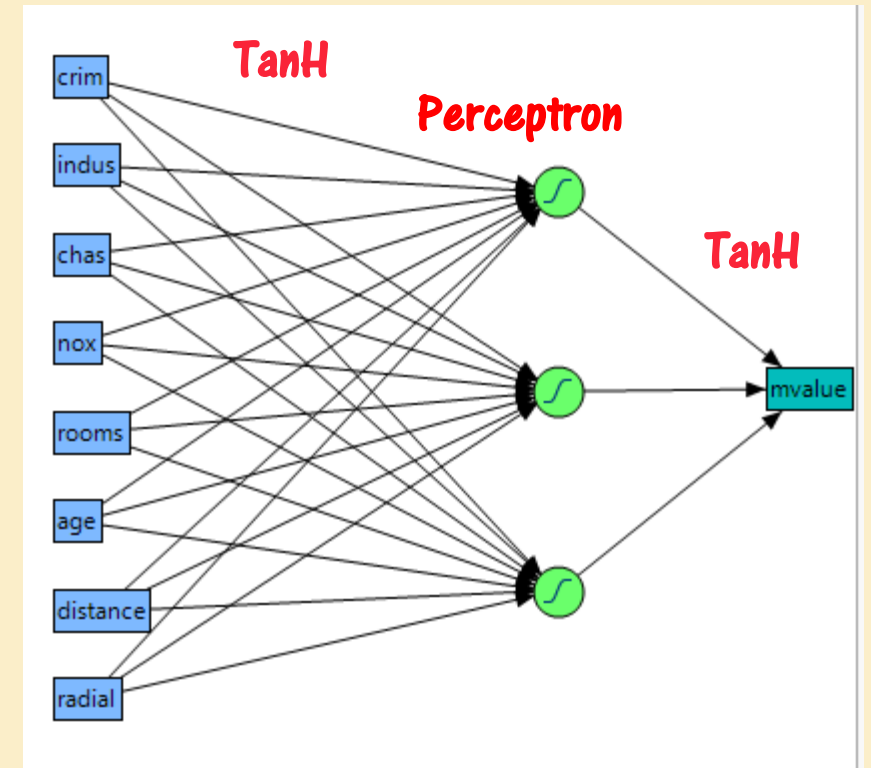
**Commands**

Missing Value Report	Number of missing values for each column
Missing Value Clustering	Hierarchical clustering of rows and columns missingness
Missing Value Snapshot	Patterns of missing values with graphical map
Multivariate Normal Imputation	Least squares prediction from the nonmissing variables in each row
Multivariate SVD Imputation	Imputation for wide problems using a singular value decomposition with the power-method adapted for missing values

- Any other better imputation method?

# JMP Neural Network Platform

- Implements a fully connected **Perceptron** (hidden nodes) with one layer.
- **Main advantage:** can efficiently model different response surfaces given enough hidden nodes and layers.
- **Main disadvantage:** results are not easily interpretable, since there are intermediate layers (**Black Box**)



## Standard JMP Edition:

- **Only TanH activation function**
- **Can fit with one hidden layer.**

**TanH** The hyperbolic tangent function is a sigmoid function. TanH transforms values to be between -1 and 1, and is the centered and scaled version of the logistic function. The hyperbolic tangent function is:

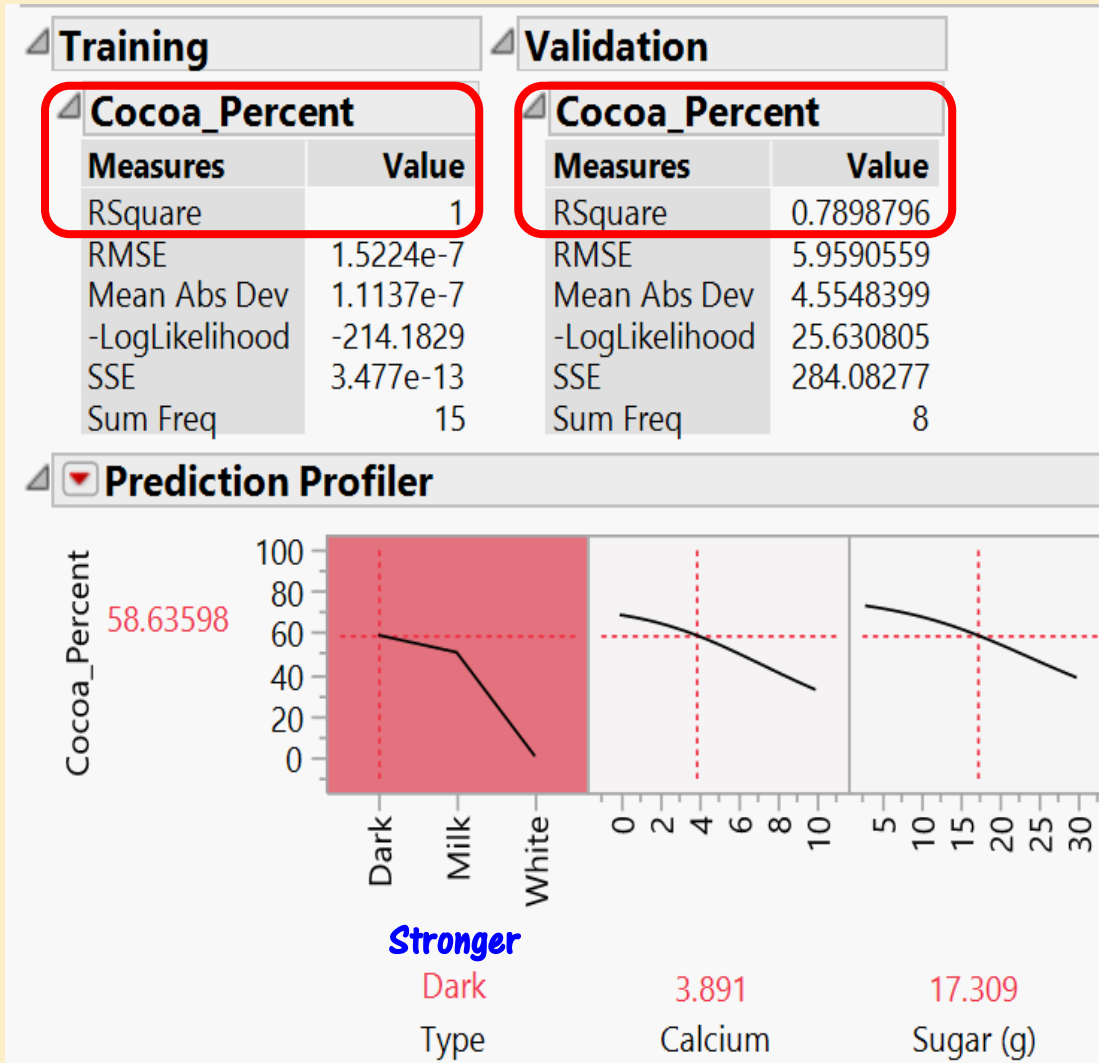
$$\frac{e^{2x} - 1}{e^{2x} + 1}$$

**More Powerful Exponential Transformation than PCA Linear Transformation**

where  $x$  is a linear combination of the  $X$  variables.

# MISSING VALUE - Neural Network

JMP 13 >> Analyze >> Predictive Modeling  
>> Neural

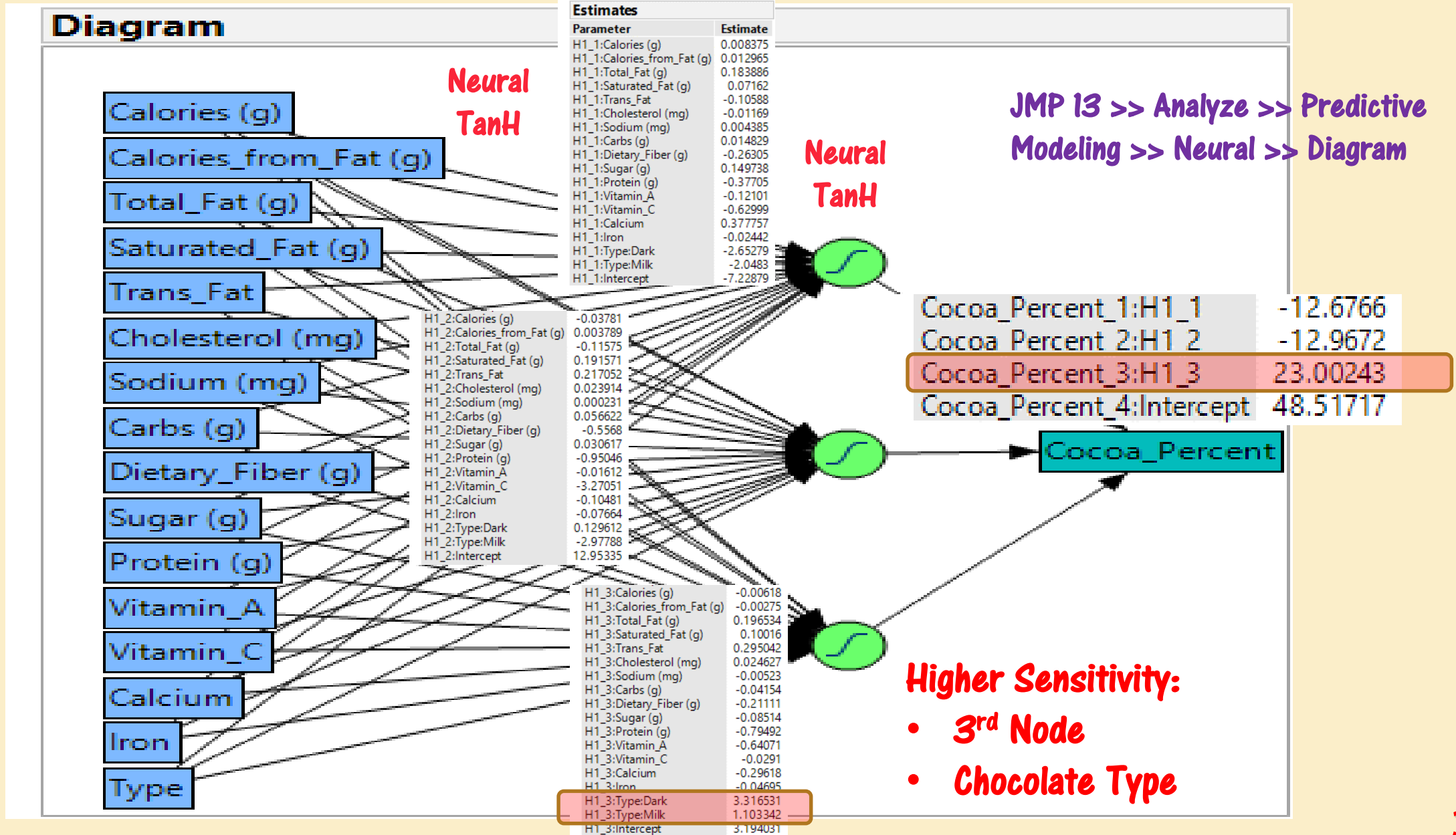


- The R-square of both Training and Validation are above 0.7.
- Though validation portion is weaker (typical Over-fit concern for Neural).
- Chocolate Type, Calcium, and Sugar as top three predictors for predicting the Cocoa%



# Neural Network- Estimates (JMP Default Setting)

AI



JMP 13 >> DOE >> Definitive Screening Design

JMP 13 >> Analyze >> Fit Y by X

# 5. DSD of Neural Setting

JMP 13 >> DOE >> Design Diagnostic >> Evaluate Design

JMP 13 >> Analyze >> Distribution

JMP 13 >> Analyze >> Fit Model

JMP 13 >> Save Script >> To Data Table or To Script Window >> Edit/Save/Run Script

- Mason C., "Optimize Neural Network Algorithm of Missing Value Imputation", submitted to 2019 ASA ENAR Spring Meeting

# Resolve Neural **Over-Fit** Concern

JMP 13 >> DOE >> Definitive Screening Design

**Objective:** optimize Neural settings to resolve over-fit by improving R-Square of both Training and Validation for **Cocoa Missing Imputation**

## JMP Neural Validation Methods:

- **Holdback:** randomly divides the original data into the training and validation (holdback portion) sets.
- **Kfold:** divides the data into K subsets. Each K set used to validate the model fit on the rest of the data, fitting a total of K models. Chose model giving the **best validation statistic**. Best for small data sets (makes efficient use of limited data)

## Four DOE Input Variables:

- Validation Method (Categorical)
- Validation Setting (Continuous) “**Nested**” under Validation Method
- Random Seed (Categorical)
- Number of Hidden Nodes (Continuous)

## Two DOE Output Responses:

- R-Square of Training Set
- R-Square of Validation Set (**More Important-Neural Over-fit**)

# Evaluate DSD of Optimizing Neural Settings

JMP 13 >> DOE >> Design Diagnostic >> Evaluate Design

Power Test of Sign (> 90%)

Correlation of Confounding (<0.3)

Uniformity of Prediction Power

14 DSD Runs

Add Four Random Corner Points



18 DSD Runs is safer on Power

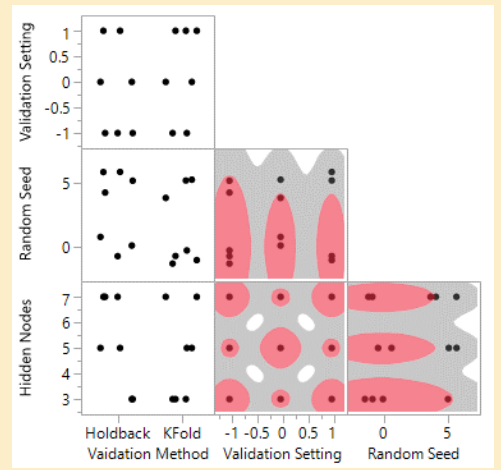
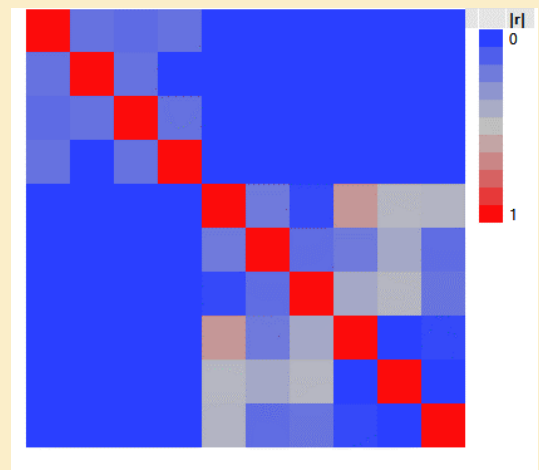
**Power Analysis**

Significance Level

Anticipated RMSE

Term	Anticipated Coefficient	Power
Intercept	1	0.913
Validation Method	1	0.89
Validation Setting	1	0.776
Random Seed	1	0.89
Hidden Nodes	1	0.783

Apply Changes to Anticipated Coefficients



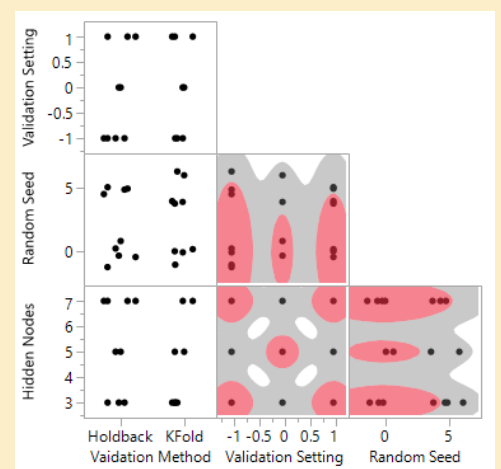
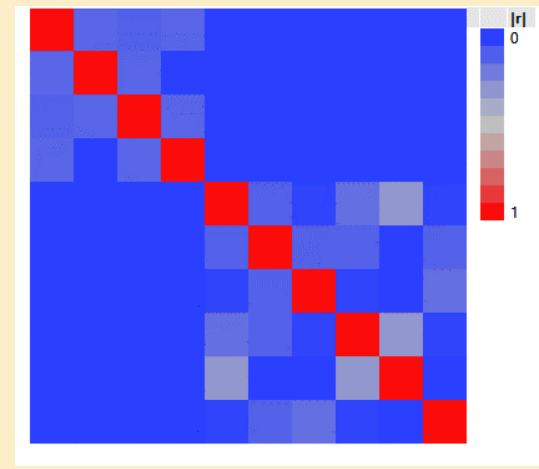
**Power Analysis**

Significance Level

Anticipated RMSE

Term	Anticipated Coefficient	Power
Intercept	1	0.975
Validation Method	1	0.97
Validation Setting	1	0.925
Random Seed	1	0.97
Hidden Nodes	1	0.925

Apply Changes to Anticipated Coefficients



# Fit Model (Nested) and Set Desirability

JMP 13 >> Analyze >> Fit Model

Pick Role Variables

Y: R-Square of Training Set, R-Square of Validation Set (optional)

Weight: optional numeric

Freq: optional numeric

By: optional

Personality: Standard Least Squares

Emphasis: Effect Screening

Fit Separately

Buttons: Help, Run, Recall, Remove

Keep dialog open:

---

Construct Model Effects

Buttons: Add, Cross, Nest, Macros

Effects List: Validation Setting[Validation Method], Validation Method, Random Seed, Hidden Nodes& RS, Validation Method\*Random Seed, Validation Method\*Hidden Nodes, Random Seed\*Hidden Nodes, Hidden Nodes\*Hidden Nodes

Degree: 2

### Construct Model Effects:

- Validation Setting is “**Nested**” under Validation Method
- Choose **Response Surface (RS)**

Maximize

R-Square of Training Set	Values	Desirability
High:	0.999	0.9819
Middle:	0.925	0.5
Low:	0.85	0.066
Importance:	1	

Maximize

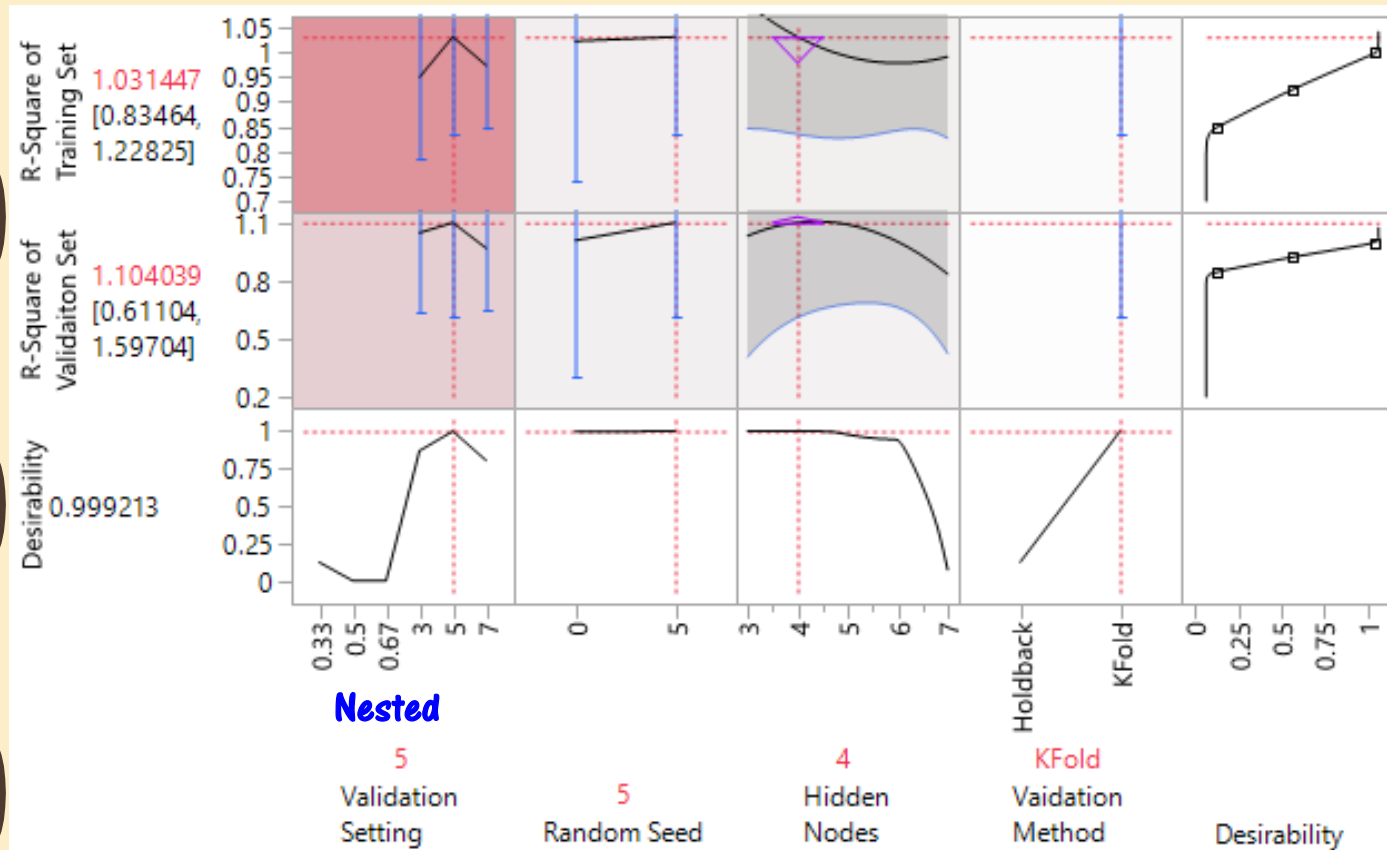
**Resolve Neural Over-Fit**

R-Square of Validation Set	Values	Desirability
High:	0.999	0.9819
Middle:	0.925	0.5
Low:	0.85	0.066
Importance:	2	



# Optimal Neural Network Setting

JMP 13 >> Analyze >> Fit Y by X



## Optimal Neural Setting:

- **Kfold** is better than Holdback (small sample size and favor validation)
- **5 Kfold** numbers (24/5 ~ 5 data for validation set)
- Use Random **Seed= 5** to improve reproducibility
- **4 Hidden Nodes** is best (Constrained by 15 input variables for one layer)
- Achieve **>99% R-Square** fit on predicting Cocoa%

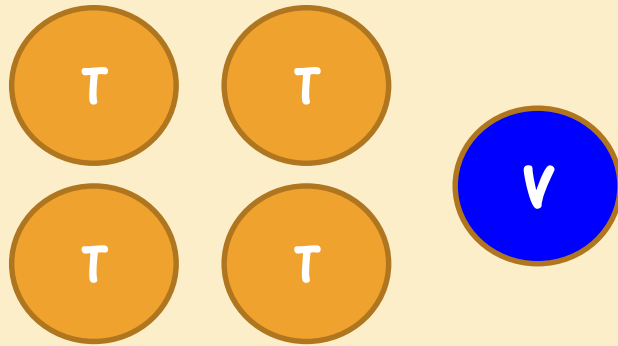
## Future Work:

- R-Square > 100%, not following Normal Distribution?
- Wider Confidence Interval: Small Validation Dataset, or Outliers?

# Understand Neural Optimization (**Future Work**)

Why Kfold over Holdback?

Holdback Portion = 0.2



Kfold **K=5**, Select the **Best**  
among **5** Models

Consider Neural Over-fit (lower Validation R-Square)

- If **K is large**, small size in each K cluster, making validation Over-Fit concern worse
- If **K is small**, losing advantage of using Kfold over Holdback
- **When total sample size is smaller, may prefer Kfold method with smaller K number**

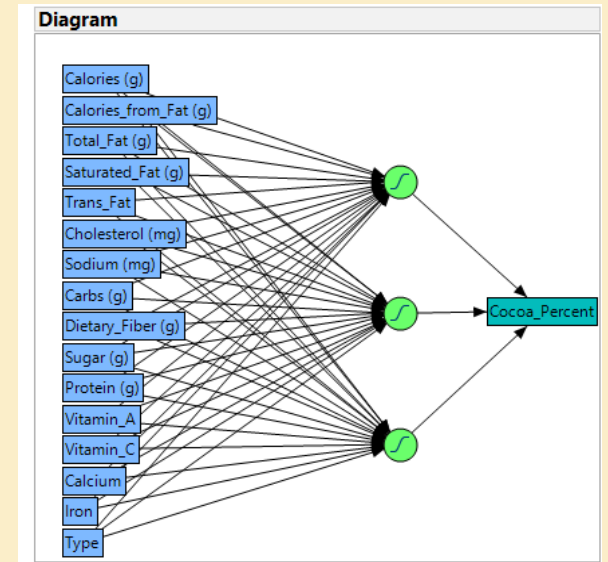
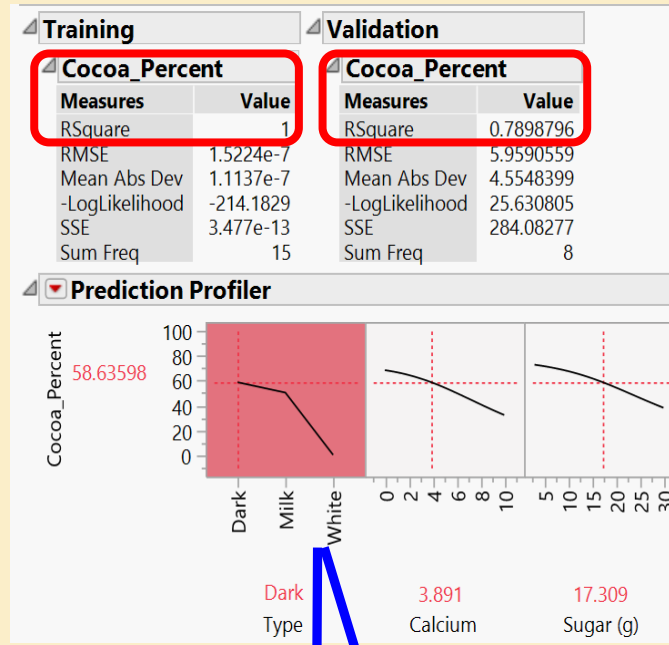
Coincidence with Four Hidden Nodes?

- The optimal Neural suggests four hidden nodes of transforming the 15 Input Nutrition Variables
- Section 2 Clustering Variables also suggests **four** clusters
- **Neural related to PCA Eigen algorithm (TanH ~ Linear)?**

# Neural Model Enhancement

JMP 13 >> Analyze >> Predictive Modeling >> Neural

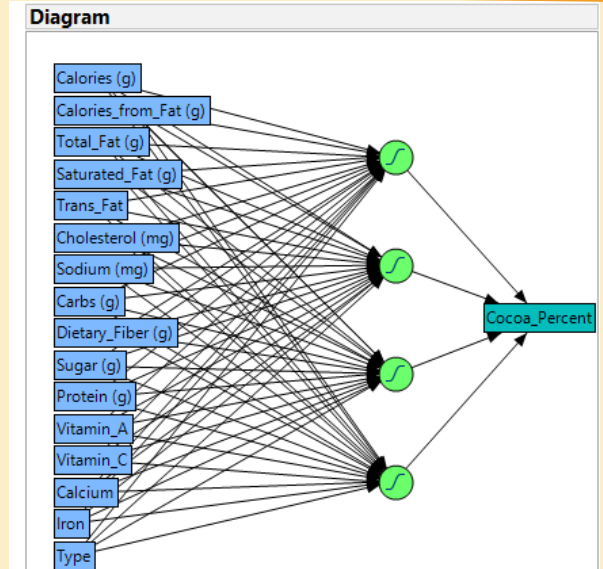
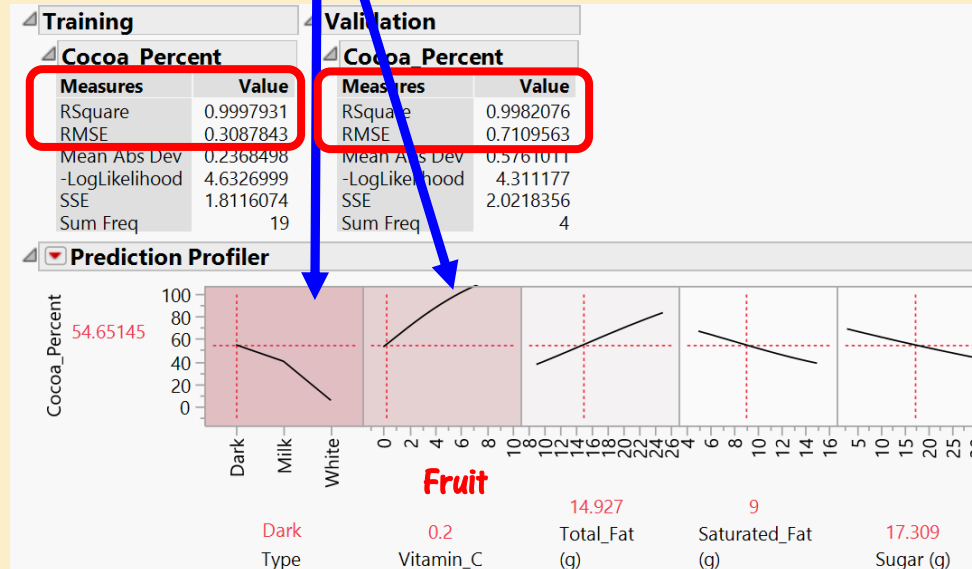
Original Neural Setting



Validation R-Square improved by 20%

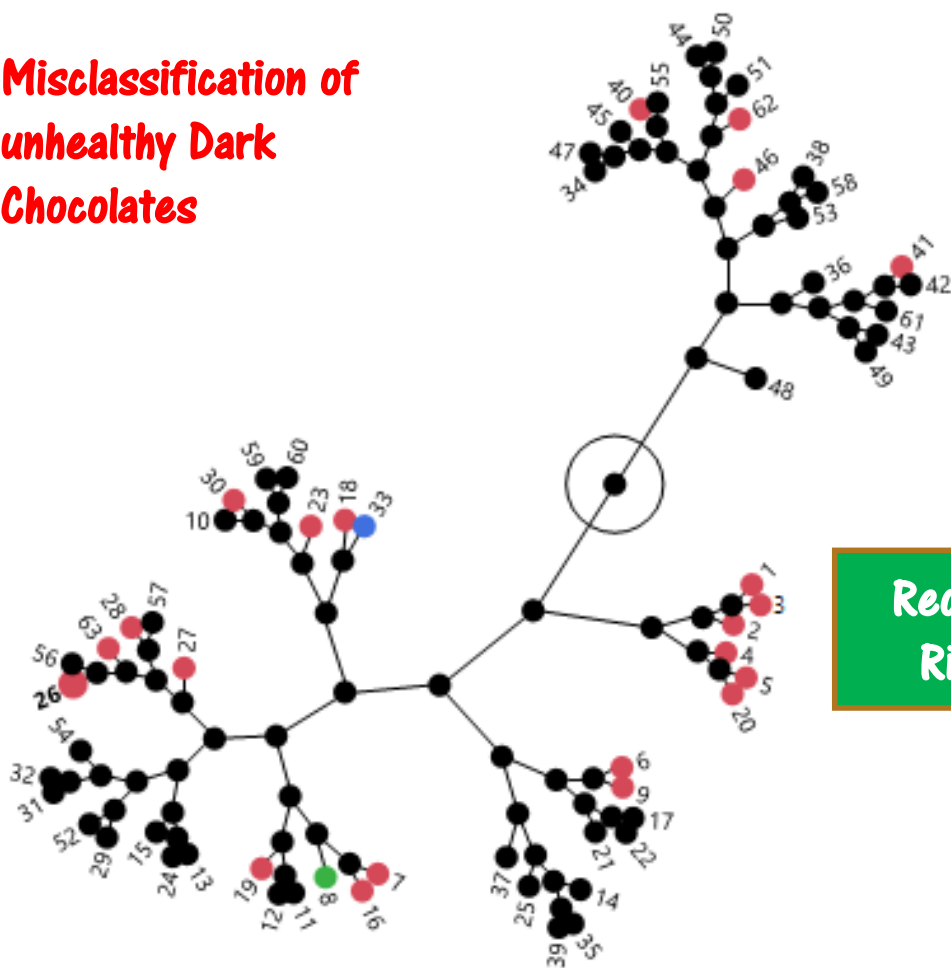
Indicating 4<sup>th</sup> Chocolate Type- Fruit Chocolate (Vitamin C)

Optimal DSD Neural Setting



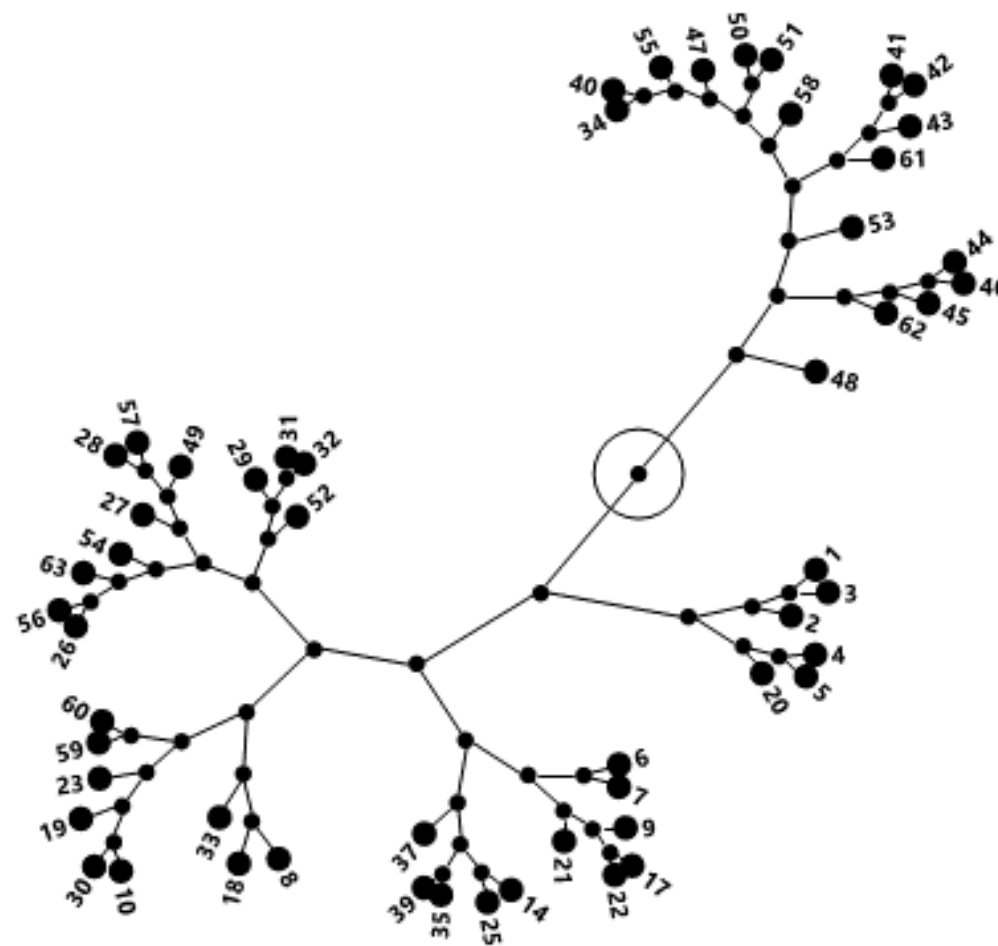
## Default: Missing Value Imputation

Misclassification of unhealthy Dark Chocolates



Reduce Risk

## Optimal Neural (Predicted Cocoa %)



# Achievements AND FUTURE RESEARCH

## Achievements:

- ✓ Adopted and Integrated “STEAMS” methodology successfully
- ✓ Learned Chocolate Products, Nutrition Anti-Oxidant Science
- ✓ Applied Multivariate Statistics, Clustering and Neural Algorithms
- ✓ Conducted DSD optimization on Resolving Neural Overfit

## Future JMP Research:

- Investigate “Fruit” Chocolate Type, Outlier Effect
- JMP Pro Partition: Bootstrap Forrest, Boosted Tree, K-Nested, Naïve Bayes
- JMP Pro Neural: Deep Learning, Hidden Layer Structure, Fitting Options
- Certify JMP Script Specialist



JAVA/Latex Advisor: Dr. Ying  
Huang

Biology/Writing Advisor: CMQ/OE  
Patrick Giuliano

Biology Advisor: Dr. I-  
Chen Chen

**Q&A**

STEAMS Advisor: ASA Dr.  
Chris Barker

STEAMS Advisor: JMP  
Chuck Boiler

**Thank YOU**

STEAMS Advisor:  
IEOM Dr. Ali Ahad,  
IEOM Dr. Don Reimer

STEAMS Advisors:  
Stanford OHS

STEAMS Advisor: ASQ Fellow Dr.  
John Flaig

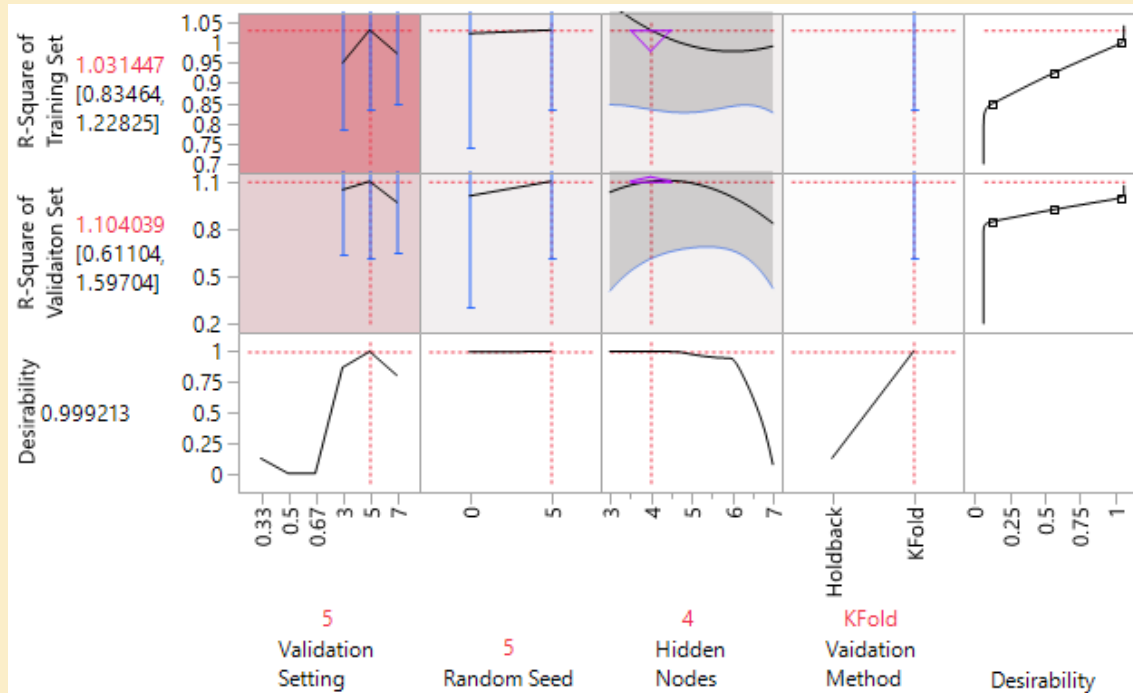
Statistics//MATH/AI Advisor: Dr.  
Charles Chen

Robotics Advisor: CQE/CRE Roland  
Jones

# Optimize Desirability Function Importance

Optimize the Number of Hidden Nodes:

- Higher R-Square of Training Set at Nodes=3
- Higher R-Square of Validation Set at Nodes=4
- Relative Importance can impact the Optimal Setting



Importance Setting		Optimization Result	Optimal Neural Settings			
Importance of Training Set	Importance of Validation Set	Overall Desirability	Validation Method	Validation Setting	Random Seed	Hidden Nodes
1	1	0.9997	KFold	5	5	3
1	1.5	0.9997	KFold	5	5	4
1	2	0.9997	KFold	5	5	4
1.5	1	0.9997	KFold	5	5	3
1.5	1.5	0.9997	KFold	5	5	3
1.5	2	0.9997	KFold	5	5	4
2	1	0.9997	KFold	5	5	3
2	1.5	0.9997	KFold	5	5	3
2	2	0.9997	KFold	5	5	3

Conduct a 2-Factor & 3-Level Full Factorial on comparing the relative importance in (1,2) range

- Set the Desirability Function Range of (0.999, 0.95, 0.9)
- In General, the optimal result shows the similar trend: 3 hidden nodes favor training set and 4 hidden nodes favor validation set

Little room for further improvement on setting the relative importance between Training Set and Validation Set

# Document Key JMP Scripts

JMP 13 >> Save Script >> To Data Table  
or To Script Window >> Edit/Save/Run  
Script

```
Partition(  
Y( :Cocoa_Percent ),  
X(  
:Type,  
:Name( "Calories (g)" ),  
:Name( "Calories_from_Fat (g)" ),  
:Name( "Total_Fat (g)" ),  
:Name( "Saturated_Fat (g)" ),  
:Trans_Fat,  
:Name( "Cholesterol (mg)" ),  
:Name( "Sodium (mg)" ),  
:Name( "Carbs (g)" ),  
:Name( "Dietary_Fiber (g)" ),  
:Name( "Sugar (g)" ),  
:Name( "Protein (g)" ),  
:Vitamin_A,  
:Vitamin_C,  
:Calcium,  
:Iron  
),  
Minimum Size Split( 3 ),  
Validation Portion( 0.6 ),  
Split History( 1 ),  
Informative Missing( 1 ),  
Column Contributions( 1 ),  
Initial Splits( :Name( "Cholesterol  
(mg)" ) >= 5 ),  
SendToReport( Dispatch( {}, "Partition",  
FrameBox, {Frame Size( 480, 56 )} ) )  
);
```

```
Neural(  
Y( :Cocoa_Percent ),  
X(  
:Name( "Calories (g)" ),  
:Name( "Calories_from_Fat (g)" ),  
:Name( "Total_Fat (g)" ),  
:Name( "Saturated_Fat (g)" ),  
:Trans_Fat,  
:Name( "Cholesterol (mg)" ),  
:Name( "Sodium (mg)" ),  
:Name( "Carbs (g)" ),  
:Name( "Dietary_Fiber (g)" ),  
:Name( "Sugar (g)" ),  
:Name( "Protein (g)" ),  
:Vitamin_A,  
:Vitamin_C,  
:Calcium,  
:Iron,  
:Type  
),  
Informative Missing( 0 ),  
Validation Method(  
"KFold", 5 ),  
Fit(  
NTanH( 4 ),  
Diagram( 1 )  
),
```

```
Fit Model(  
Y( :Name( "R-Square of Training Set" ),  
:Name( "R-Square of Validation Set" ) ),  
Effects(  
:Validation Setting[:Validation Method],  
:Validation Method,  
:Random Seed,  
:Hidden Nodes & RS,  
:Validation Method * :Random Seed,  
:Validation Method * :Hidden Nodes,  
:Random Seed * :Hidden Nodes,  
:Hidden Nodes * :Hidden Nodes  
),  
Personality( "Standard Least Squares" ),  
Emphasis( "Effect Screening" ),  
:Name( "R-Square of Training Set" ) << {Summary of Fit( 0 ),  
Analysis of Variance( 0 ), Lack of Fit( 0 ), Sorted  
Estimates( 0 ),  
Plot Actual by Predicted( 1 ), Plot Regression( 0 ),  
Plot Residual by Predicted( 1 ), Plot Studentized Residuals(  
1 ),  
Plot Effect Leverage( 0 ), Box Cox Y Transformation( 1 )},  
:Name( "R-Square of Validation Set" ) << {Summary of Fit( 0  
) ,  
Analysis of Variance( 0 ), Lack of Fit( 0 ), Sorted  
Estimates( 0 ),  
Plot Actual by Predicted( 1 ), Plot Regression( 0 ),  
Plot Residual by Predicted( 1 ), Plot Studentized Residuals(  
1 ),  
Plot Effect Leverage( 0 ), Box Cox Y Transformation( 42 )}  
),
```