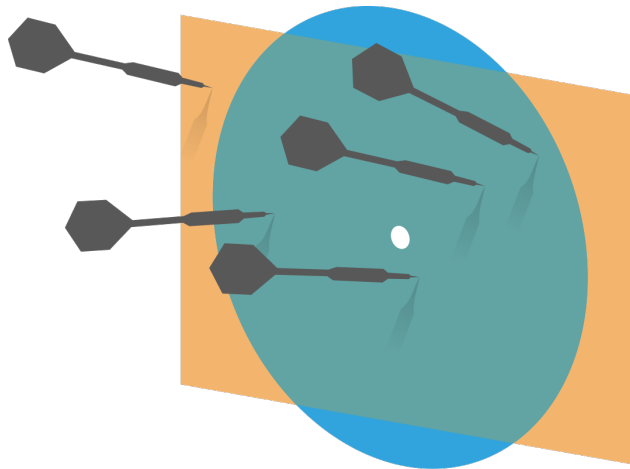# Probability and Statistics for Computer Science ↗



Credit: wikipedia

"Unsupervised learning is arguably more typical of human and animal learning…"--- Kelvin Murphy, former professor at UBC

Hongye Liu, Teaching Assistant Prof, CS361, UIUC, 12.03.2020

# Last time

* Curse of dimensions

* Unsupervised learning

* Clustering

# Objectives

* Application of Clustering

Cluster center Histogram

* Markov chain (1)

conditional prob.
coming back is Matrix

# Q. Is k-means clustering deterministic?

A. Yes
B. No

# K-means clustering example: Portugal consumers

✳ The dataset consists of the annual grocery spending of 440 customers

✳ Each customer's spending is recorded in 6 features:
  ✳ fresh food, milk, grocery, frozen, detergents/paper, delicatessen

✳ Each customer is labeled by: 6 labels in total
  ✳ Channel (Channel 1 & 2) (Horeca 298, Retail 142)
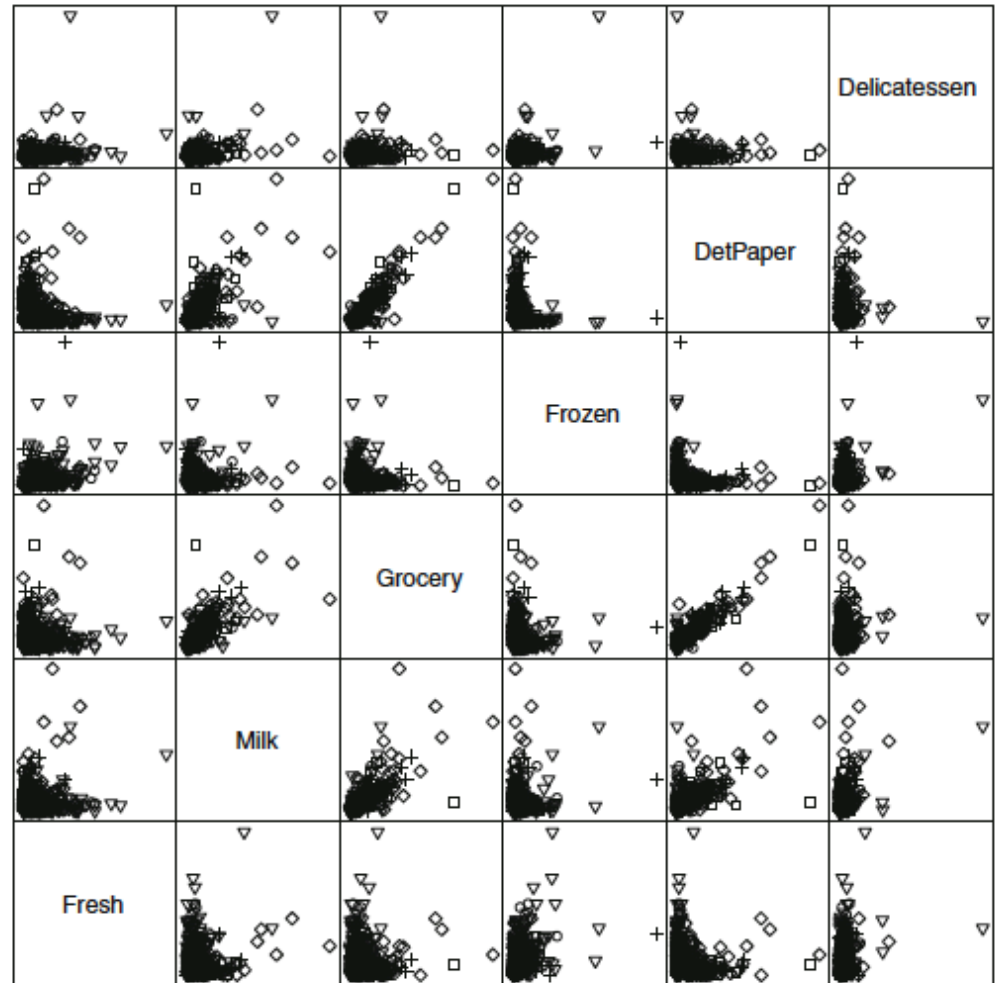  ✳ Region (Region 1, 2 &3) (Lisbon 77, Oporto 47, Other 316)

# Lisbon, Portugal

# Oporto, Portugal

# Visualization of the data

* Visualize the data with scatter plots

* We do see that some features are correlated.

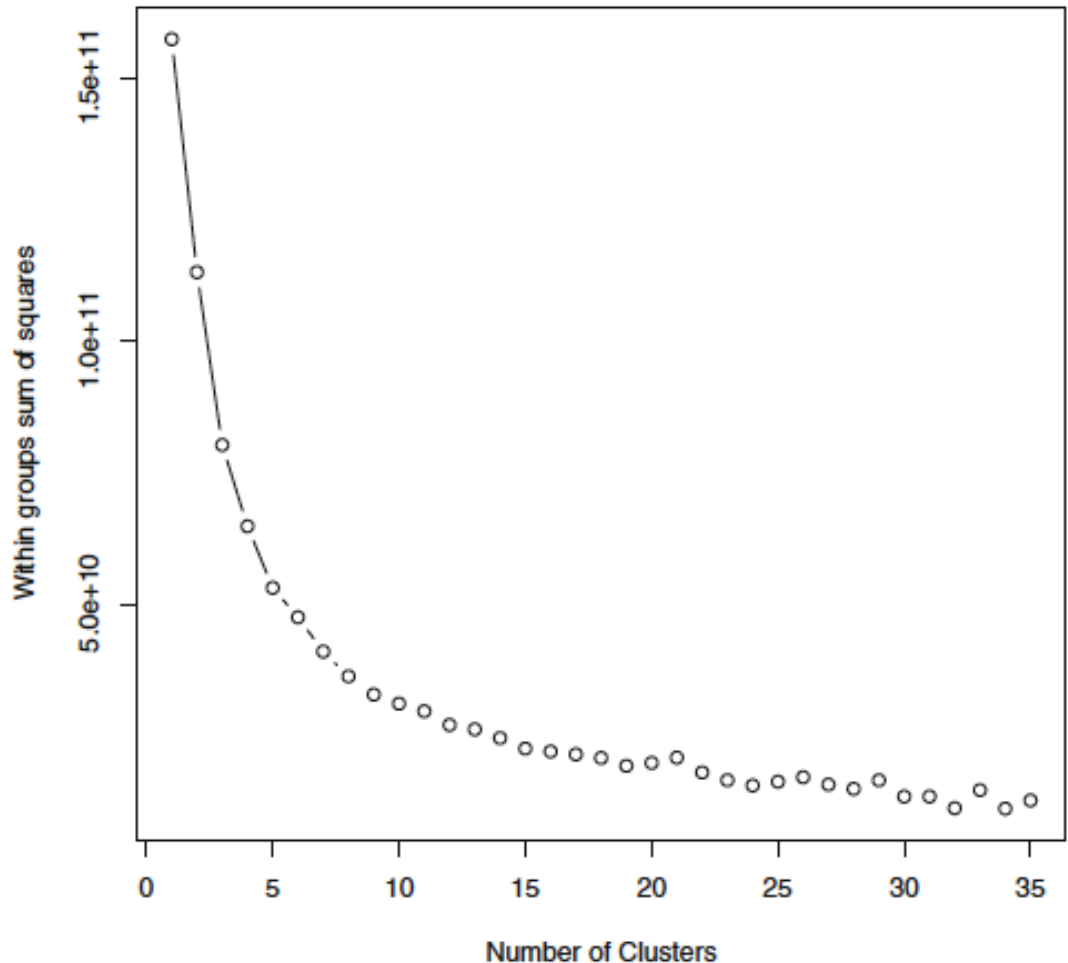* But overall we do not see significant structure or groups in the data.



Scatter Plot Matrix
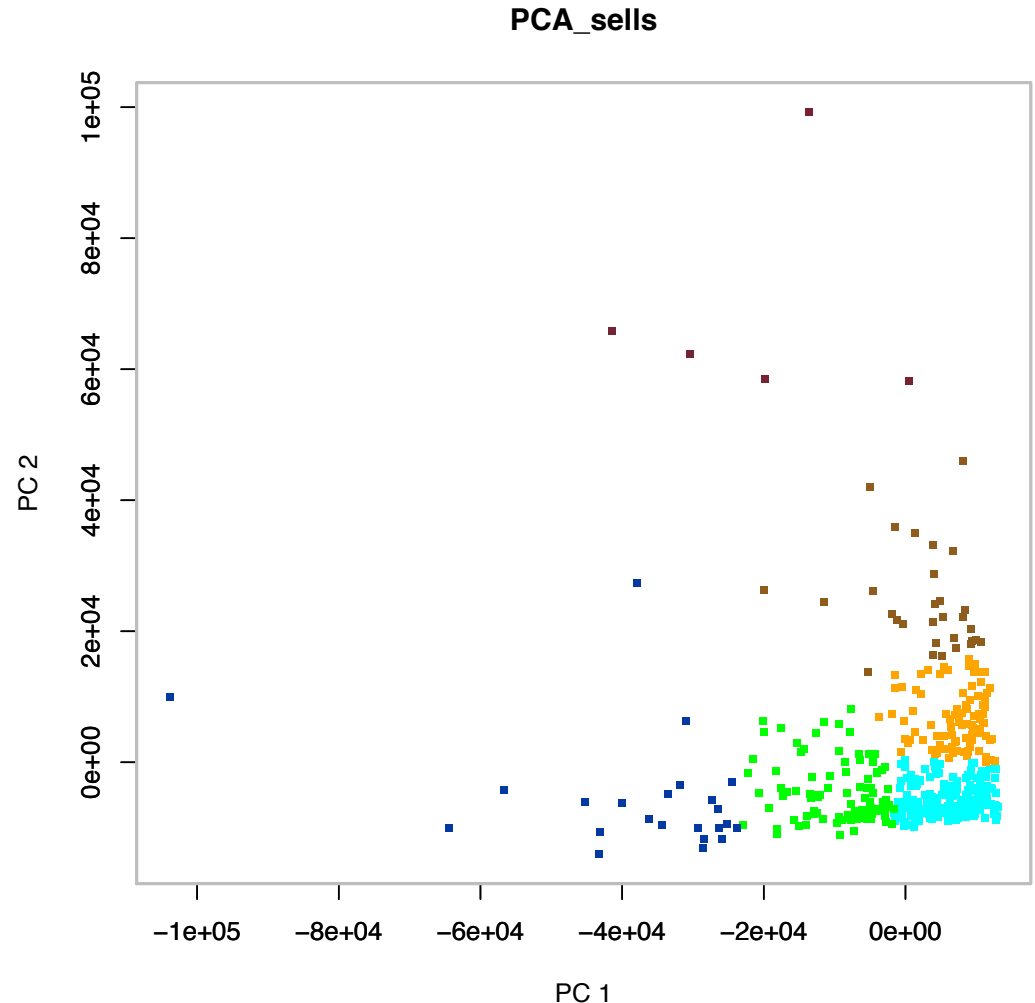
# Do kmeans and choose k through the cost function

It's good to pick a **k** around the knee:
I choose 6 for it matches the number of labels
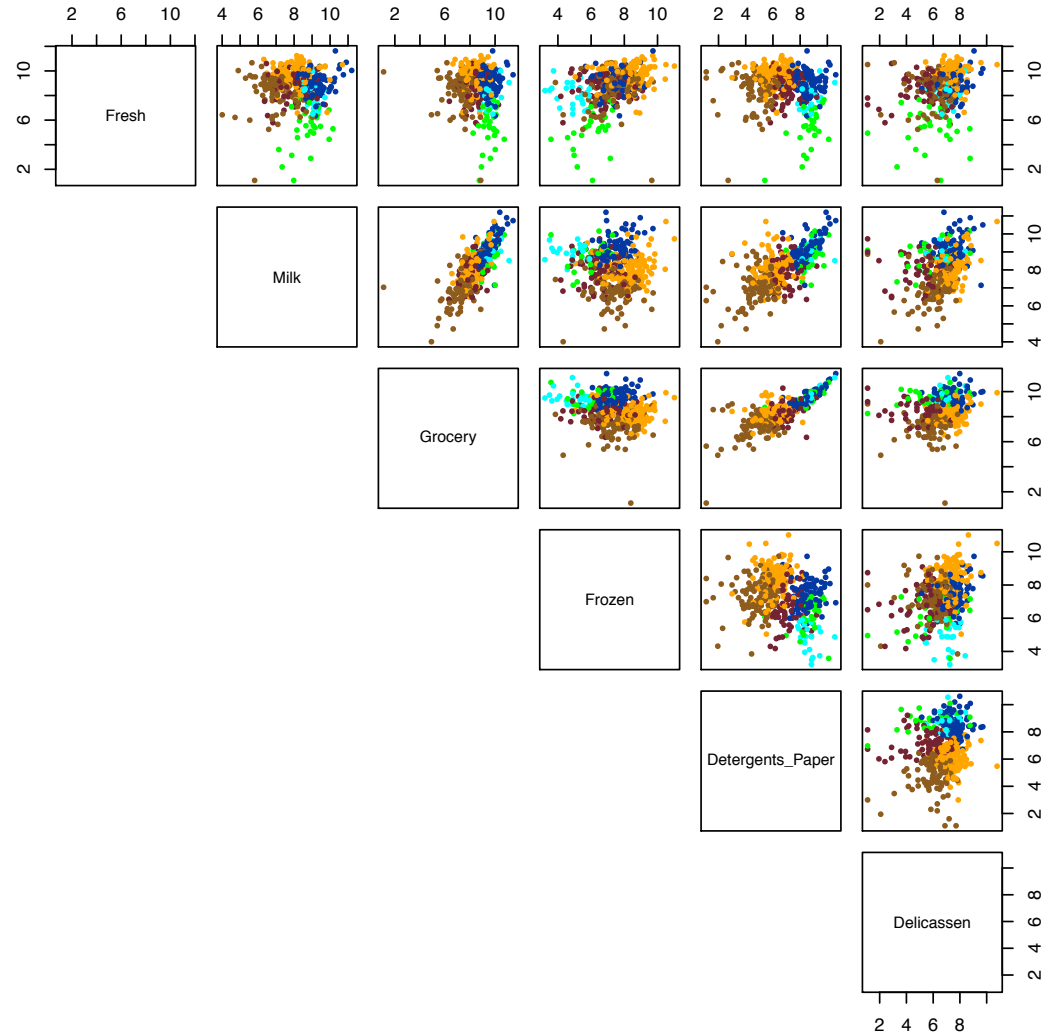
# Visualization of the data (PCA)

✳ PCA does show some separation. **Colors are the clusters**

✳ Data points show large range of dynamics!

440 customers

each dot is one customer

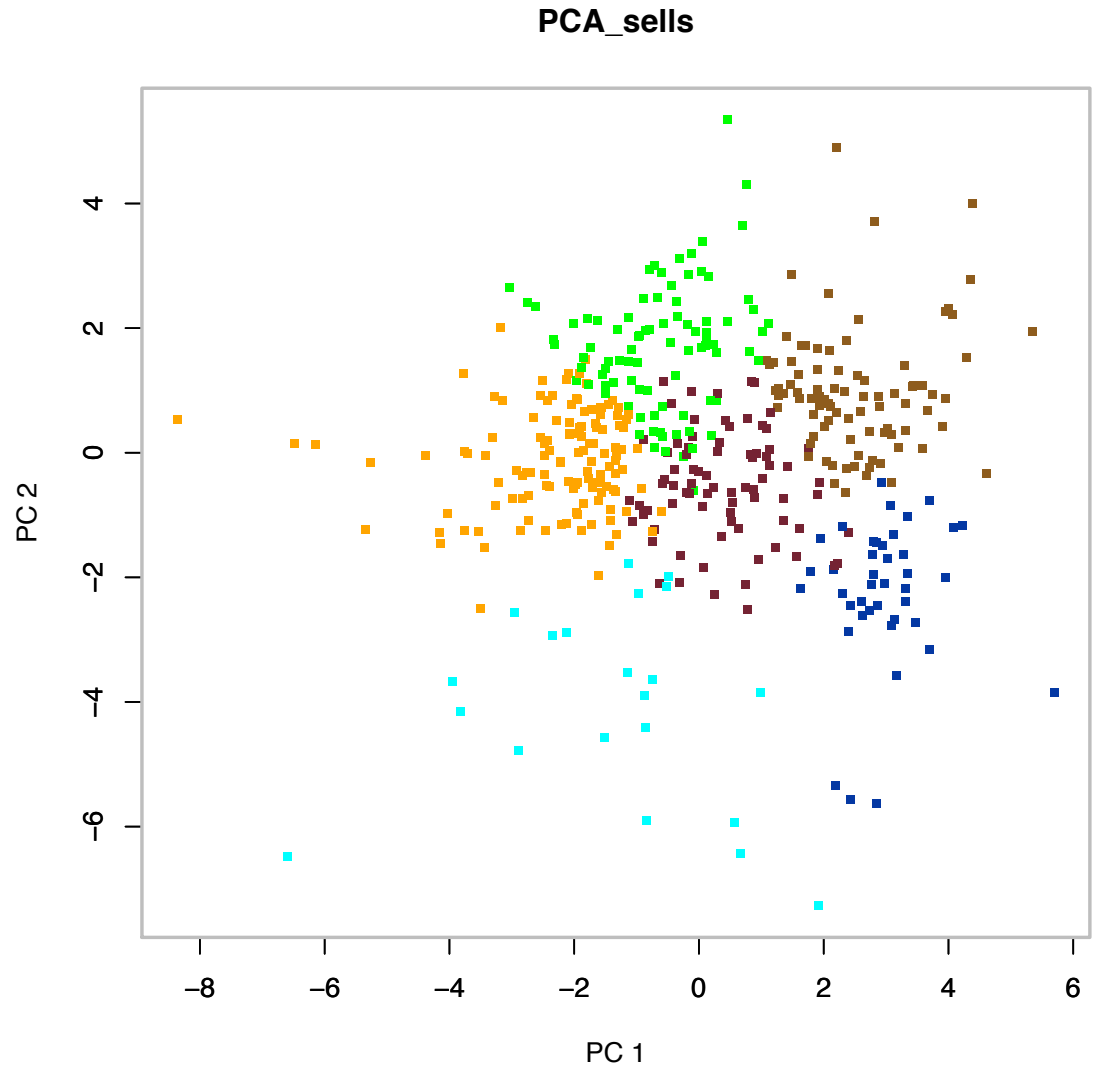

PCA_sells

# Do log transform of the data

✳ Log transform the data

✳ Do scatter plot matrix after the log transform

✳ Do the kmeans and color the clusters identified by k-means
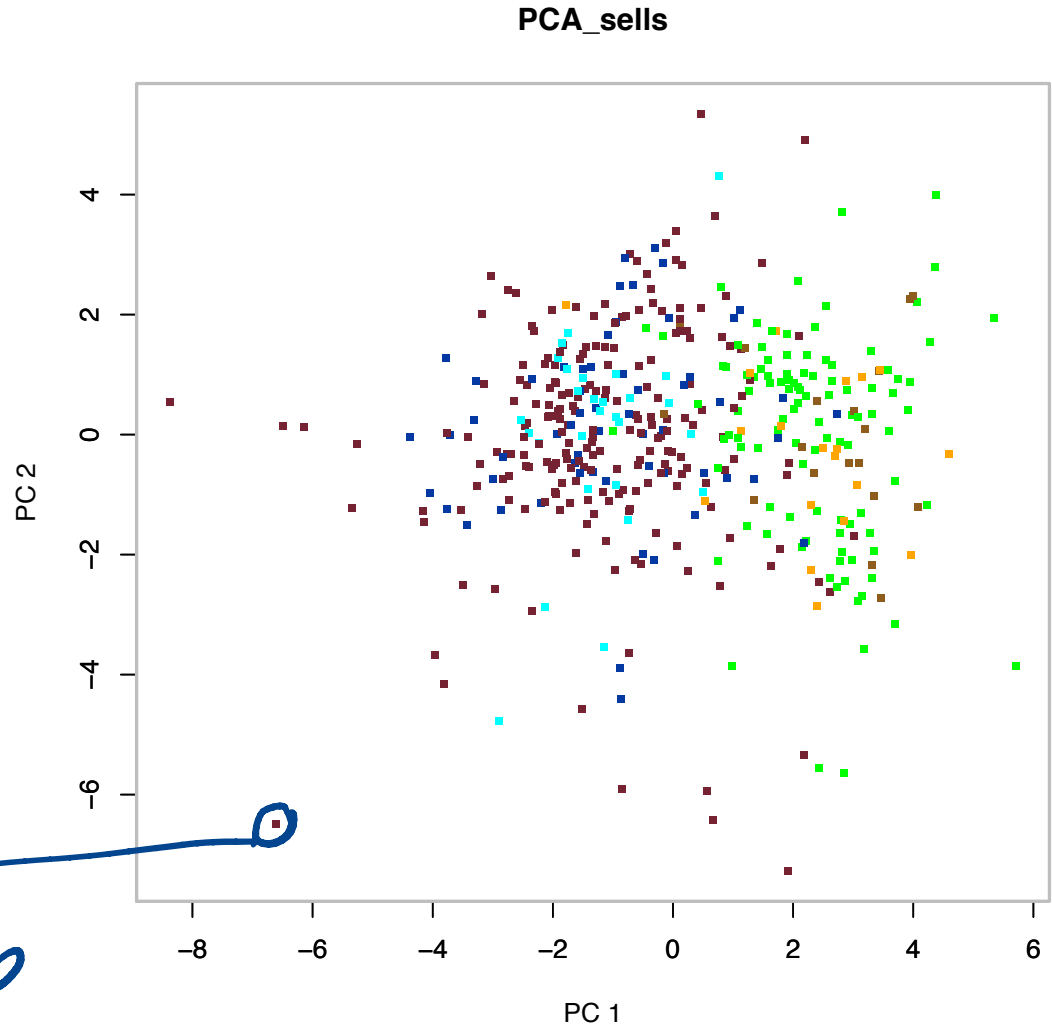
# PCA after log transformation: Clusters

Colors show the **clusters** identified by k-means



PCA_sells

# PCA after log transformation
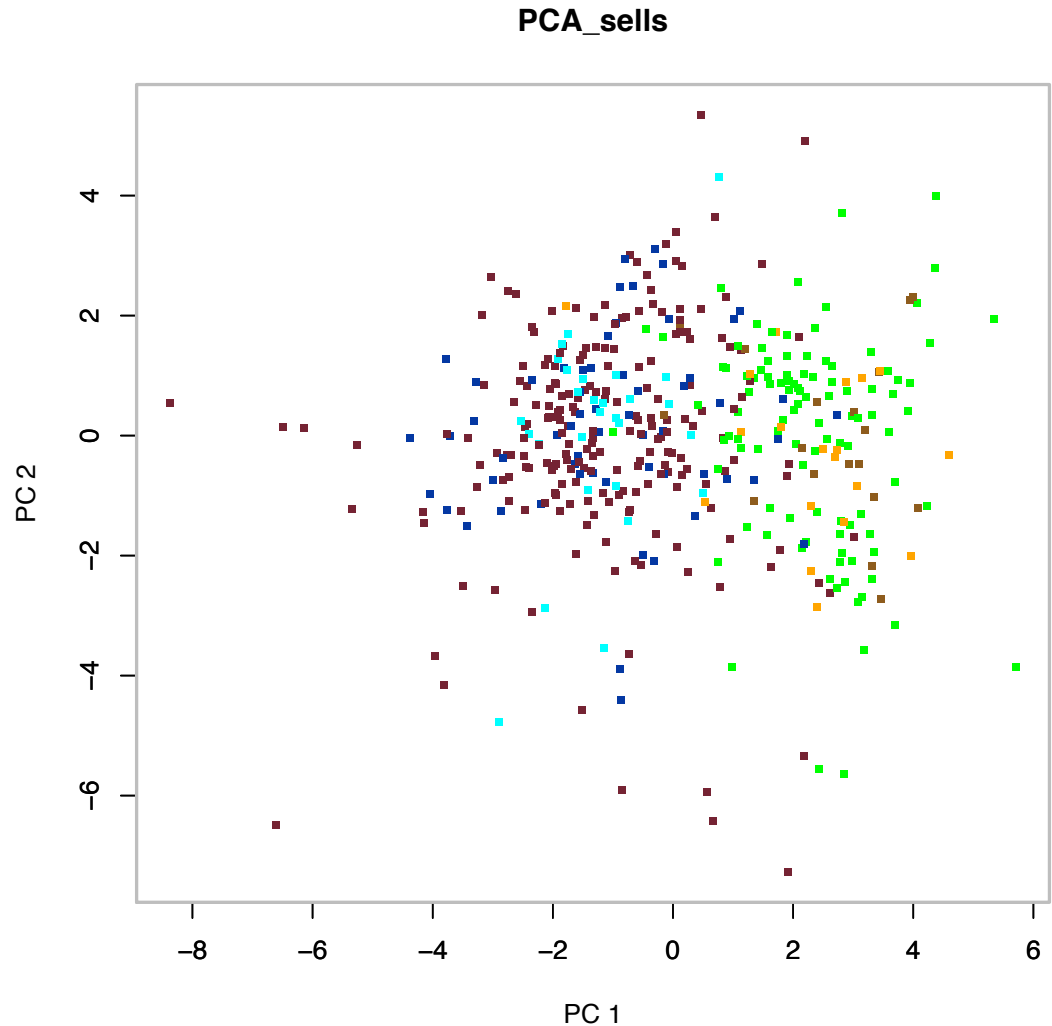
Colors show the **Channel-region labels**

What does this tell us?

Customer

440



PCA_sells

# PCA after log transformation

Colors show the **Channel-region labels**

Channels differ a lot



PCA_sells

# Cluster center histogram of the Portugal grocery spending data

✳ For each channel/ region, we make a histogram of customers that map to each of the **6 cluster centers**.

✳ **What do you see?**

Channel1: Horeca
Channel2: Retail

Region1: Lisbon
Region2: Oporto
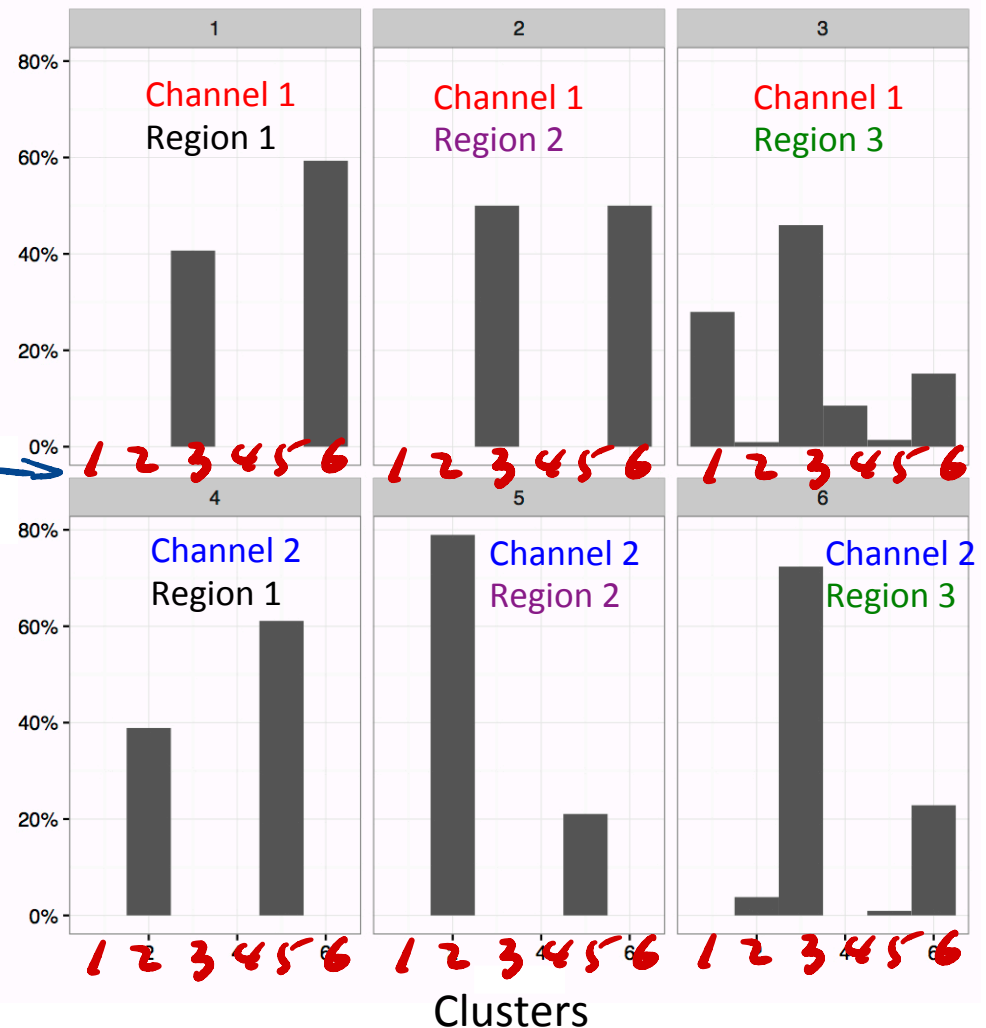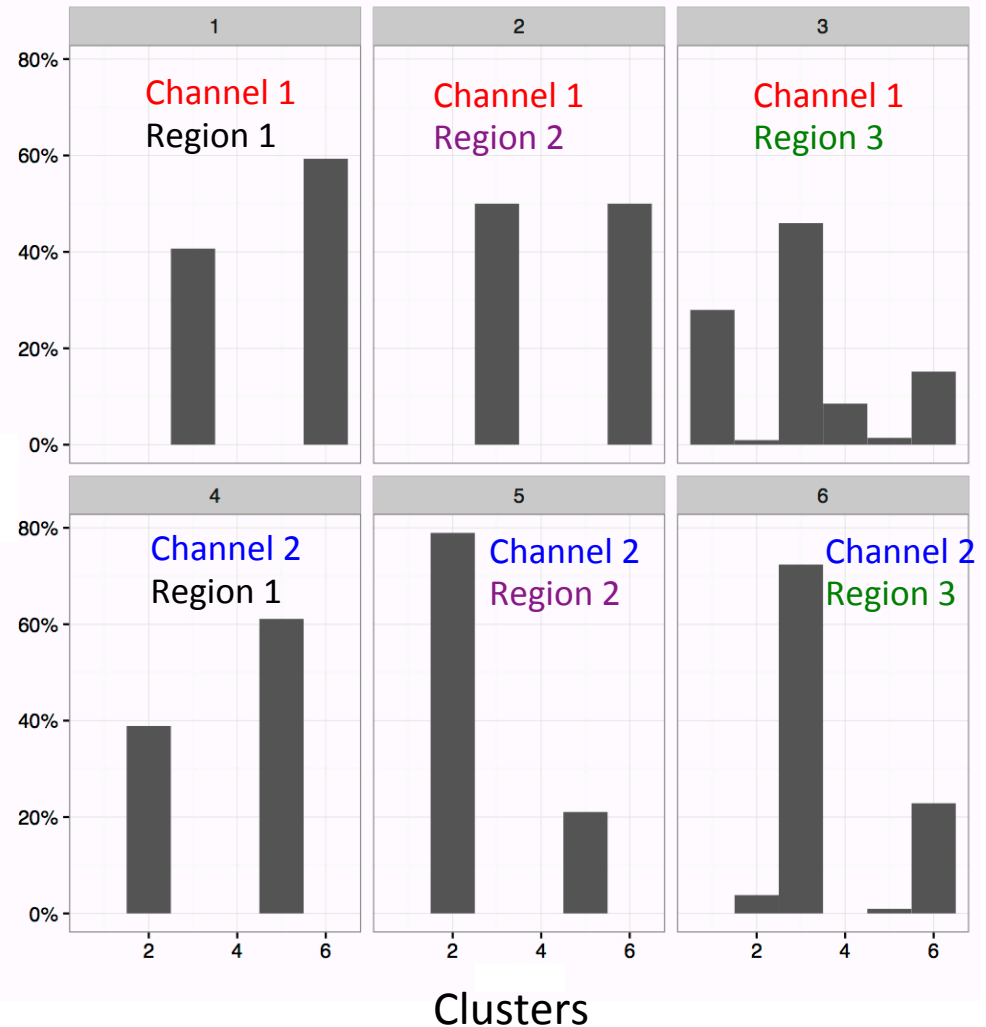Region3: Other



Clusters

# Cluster center histogram of the Portugal grocery spending data

* For each channel/ region, we make a histogram of customers that map to each of the 6 cluster centers.

* **Channels are significantly different!**

* **Region 3 is special**

* **Is it enough to plot the percentage?**

# Cluster center histogram of the Portugal grocery spending data
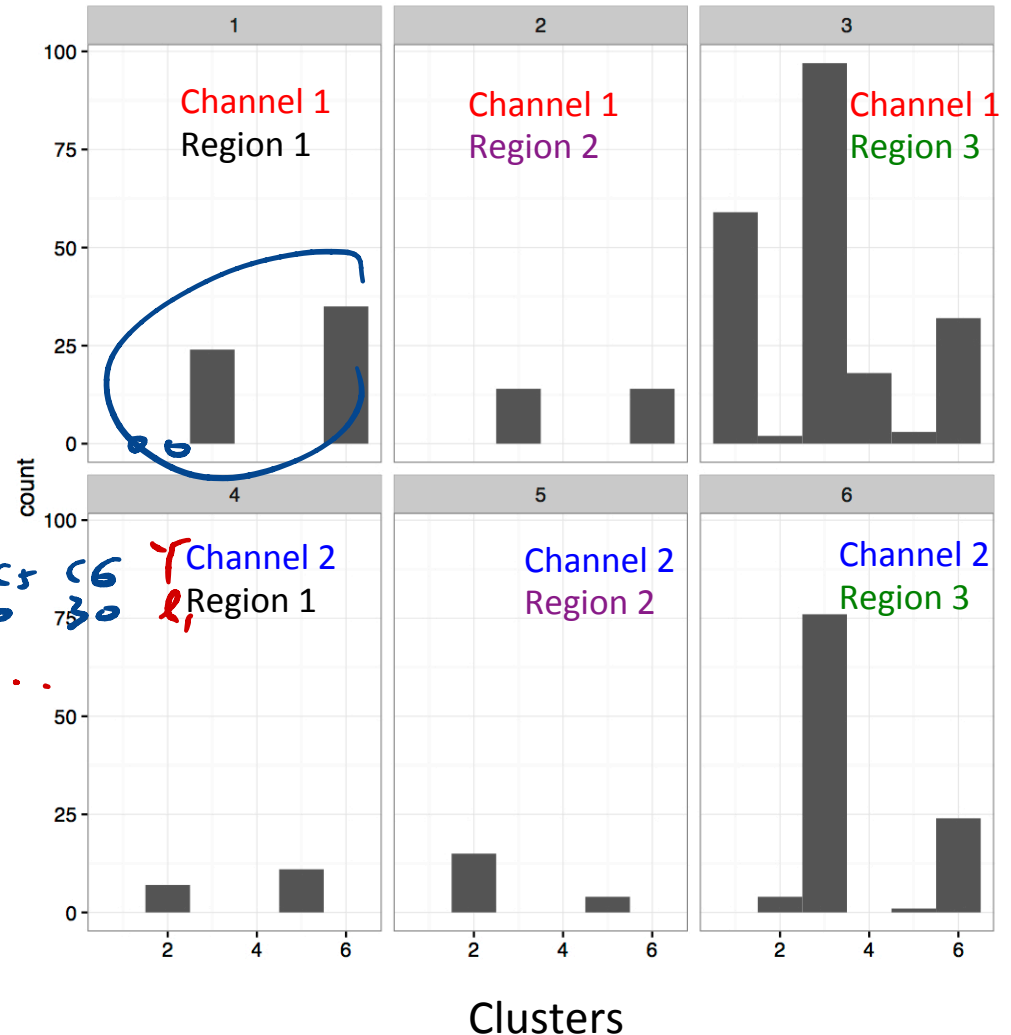
* For each channel/region, we make a histogram of customers that map to each of the 6 cluster centers.

* **Channels are significantly different!**

* **Region 3 is special**

* **Count matters depending on the purpose**

## Q. What can we do with cluster center histograms?

A. investigate the feature patterns of data groups

B. Classify new data with the cluster center histograms.

C. Both A and B.

# Markov chain

In a class, students are either up-to-date or behind regarding progress. If a student is up-to-date, the student has 0.8 probability remaining up-to-date, if a student is behind, the student has 0.6 probability becoming up-to-date. Suppose the course is so long that it runs life long, what is the probability any student eventually gets up-to-date?

A   25%

B   50%

C   25%     75%

D   9.%

# Markov Chain

* Motivation
* Definition of Markov model
* Graph representation – Markov chain
* Transition probability matrix
* The stationary Markov chain
* The pageRank algorithm

# An example of dependent events in a sequence

I had a glass of wine with my grilled _____
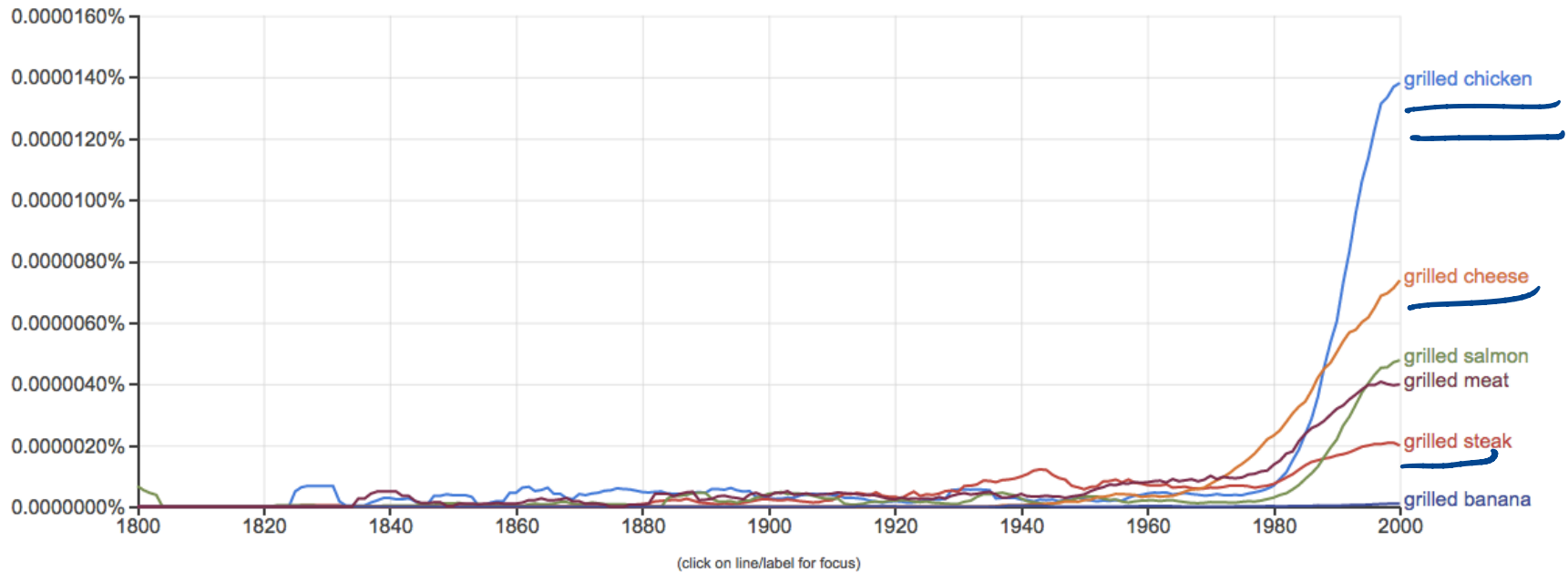
meat ✓
fish ✓
cheese ✓

wood

# An example of dependent events in a sequence

# An example of dependent events in a sequence

# Markov chain

✳ Markov chain is a process in which outcome of any trial in a sequence is **conditioned by the outcome of the trial immediately preceding, but not by earlier ones**. $P(X_n)$ $X_t$

✳ Such dependence is called **chain dependence** $P(X_{n-1})$ $X_{(t-1)}$
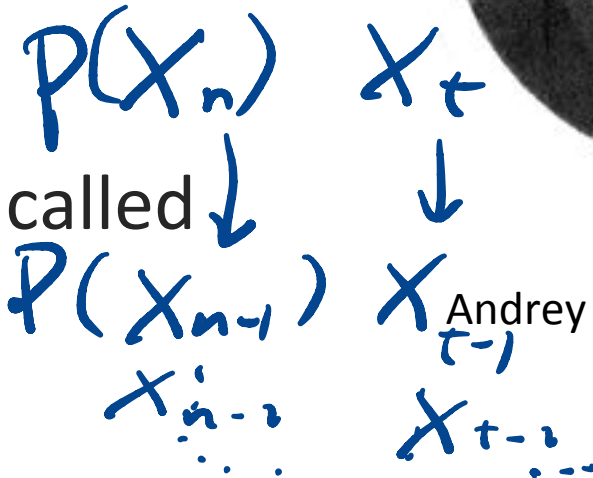
$X_{n-2}$ $X_{t-2}$

Andrey Markov (1856-1922)

# Markov chain in terms of probability

* Let $X_0$, $X_1$,… be a sequence of discrete finite-valued random variables

* The sequence is a Markov chain if the probability distribution $X_t$ only depends on the distribution of the immediately preceding random variable $X_{t-1}$

$$P(X_t|X_0..., X_{t-1}) = P(X_t|X_{t-1})$$

* If the conditional probabilities (transition probabilities) do **NOT change with time**, it's called **constant Markov chain**.

$$P(X_t|X_{t-1}) = P(X_{t-1}|X_{t-2}) = ... = P(X_1|X_0)$$

* Toss a fair coin until you see two heads in a row and then stop, what is the probability of stopping after exactly **n** flips?

✳ ✳ ✳ HH

n

Random variable

$P(n = n_0) = ?$

Geometric

T T T... H

allele    State diagram

1 -> **Start or just had tail/restart**
2 -> **had one head after start/restart**
3 -> **2heads in a row/Stop**

directed graph



T 1/2

H 1/2          H 1/2

T

1      2      3

1/2

$N =$ ① ② ③ ④ ⑤ ⑥

Trials    T   T   H   T   H   H

$X_N =$   $X_1$   $X_2$   $X_3$   $X_4$   $X_5$   $X_6$

State    1   1   2   1   2   3

Markov property:

Given the current state
the past doesn't matter

$$P_{ij} = P(X_{n+1} = j \mid X_n = i)$$

$$= P(X_{n+1} = j \mid X_n = i, \boxed{X_{n-1} = ? \quad \cdots \quad X_0 = ?})$$

This part can be
any !!

✳ Let $p_n$ be the probability of stopping after **n** flips

$$p_1 = 0 \quad p_2 = \boxed{1/4} \quad p_3 = \boxed{1/8} \quad p_4 = \boxed{1/8} \quad \cdots$$

*(handwritten annotations in red:)*

Under $p_2$: HHH

Under $p_3$: ✳ H H, T H H

Under $p_4$: ✳ ✳ HH, T T HH, H T HH

$$P(n = n_0) =$$

$n \uparrow$

✳ Let $p_n$ be the probability of stopping after **n** flips

$$p_1 = 0 \quad p_2 = 1/4 \quad p_3 = 1/8 \quad p_4 = 1/8 \quad \text{...}$$

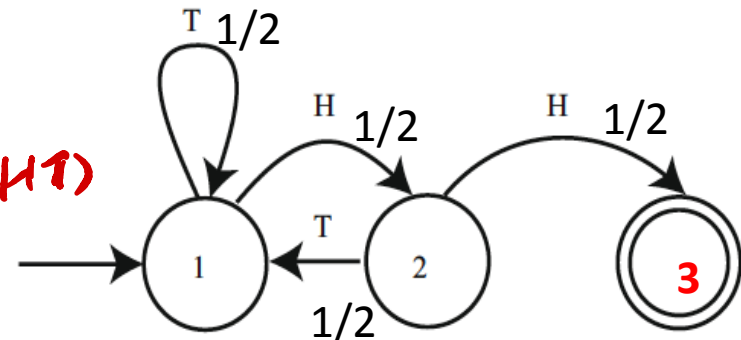✳ If $n > 2$, there are two ways the sequence starts

 ✳ Toss T and finish in n-1 tosses

 ✳ Or toss HT and finish in n-2 tosses

*(handwritten annotations: T, n-1, n; HT, n-2)*

✳ So we can derive a recurrence relation

$$p_n = \frac{1}{2}p_{n-1} + \frac{1}{4}p_{n-2}$$

*(handwritten: P(n-1|T), P(n-2|HT))*

P(T)     P(HT)

*(diagram: states 1, 2, 3 with transitions — T 1/2 self-loop, H 1/2, T 1/2, H 1/2)*

# Transition probability btw states

States

|   | ① | ② | ③ |
|---|---|---|---|
| ① | $\frac{1}{2}$ | $\frac{1}{2}$ | 0 |
| ② | $\frac{1}{2}$ | 0 | $\frac{1}{2}$ |
| ③ | 0 | 0 | 1 |

|   | ① | ② | ③ |
|---|---|---|---|
| ① | ? $\frac{3}{4}$ | $\frac{1}{4}$ | 0 |
| ② | $\frac{3}{4}$ | 0 | $\frac{1}{4}$ |
| ③ | 0 | 0 | 1 |

$P(②|①) = \frac{1}{4}$

Fair coin tosses

T 1/2

H 1/2     H 1/2

T

1     2     3

1/2

Biased coin tosses

T $\frac{3}{4}$

H $\frac{1}{4}$     H $\frac{1}{4}$

T

1     2     3

$\frac{3}{4}$

# Transition probability matrix: weather model

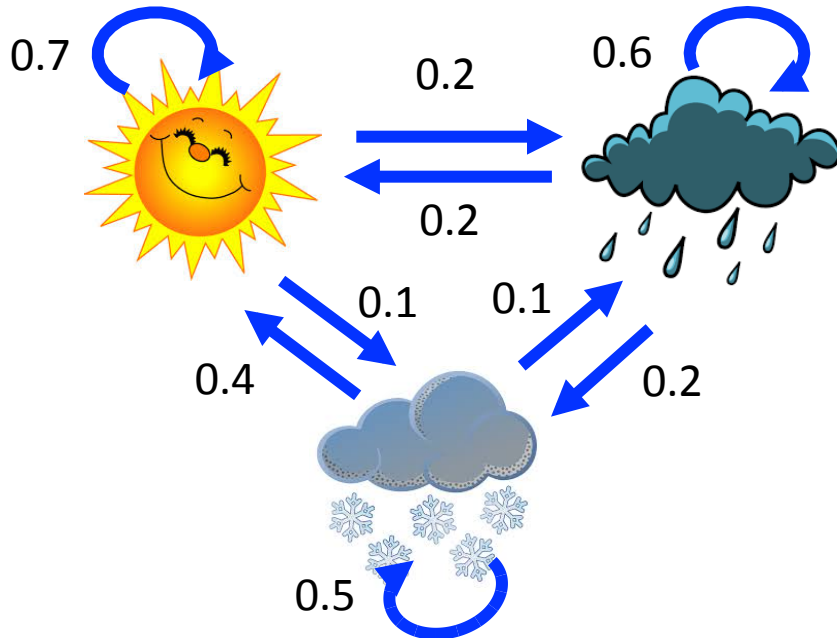✳ Let's model daily weather as one of the three states (Sunny, Rainy, and Snowy) with Markov chain that has the transition probabilities as shown here.
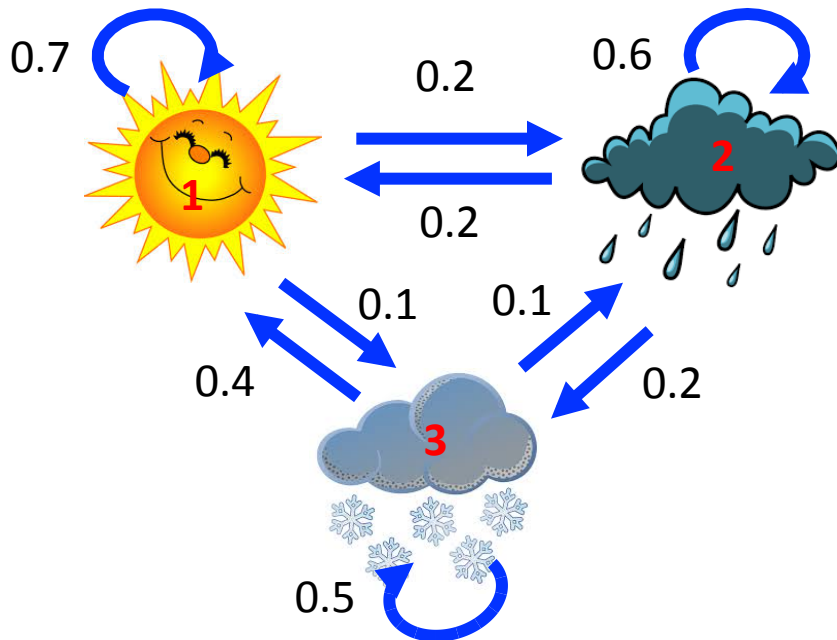
# Transition probability matrix: weather model

✳ Let's model daily weather as one of the three states (Sunny, Rainy, and Snowy) with Markov chain that has the transition probabilities as shown here.



**i**, the current state at time point t
**j**, the next state at time point t+1

$$P = \begin{array}{ccc} & \text{Sunny} & \text{Rainy} & \text{Snowy} \end{array}$$

$$P = \begin{matrix} Su \\ R \\ Sn \end{matrix} \begin{bmatrix} 0.7 & 0.2 & 0.1 \\ 0.2 & 0.6 & 0.2 \\ 0.4 & 0.1 & 0.5 \end{bmatrix} \begin{matrix} \text{Sunny} \\ \text{Rainy} \\ \text{Snowy} \end{matrix}$$

The transition probability matrix

# Q: The transition probabilities for a node sum to 1

## A. Yes.

## B. No.

Only the row sum is 1, that is: the probabilities associated with outgoing arrows sum to 1.

$P$ is a transition prob. matrix

$$\pi_0 = [ \overset{\text{Sunny}}{0} \quad \overset{\text{Rainy}}{1} \quad \overset{\text{Snowy}}{0} ] \qquad \tau$$

$$P = \begin{bmatrix} 0.7 & 0.2 & 0.1 \\ 0.2 & 0.6 & 0.2 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$$

$P(\text{Snowy})$ for next day? $\qquad t+1$

$$\pi_1 = \pi_0 P$$

$$= [0 \ 1 \ 0] \begin{bmatrix} & & \end{bmatrix} = [ \qquad \overset{\text{Sunny Rainy Snowy}}{\underset{\uparrow}{0.2}} ]$$

# Additional References

✳ Robert V. Hogg, Elliot A. Tanis and Dale L. Zimmerman. "Probability and Statistical Inference"

✳ Kelvin Murphy, "Machine learning, A Probabilistic perspective"

*See You!*