# Probability and Statistics for Computer Science
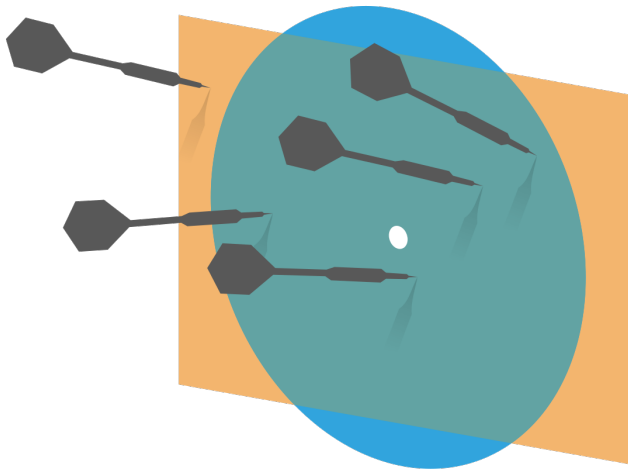


Credit: wikipedia

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

Hongye Liu, Teaching Assistant Prof, CS361, UIUC, 04.29.2021

# Last time

* Curse of dimensions

* Unsupervised learning

* Clustering

# Objectives

\* Application of Clustering

Cluster Center Histogram

\* Markov Chain (1)

Conditional probability

coming back in 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)

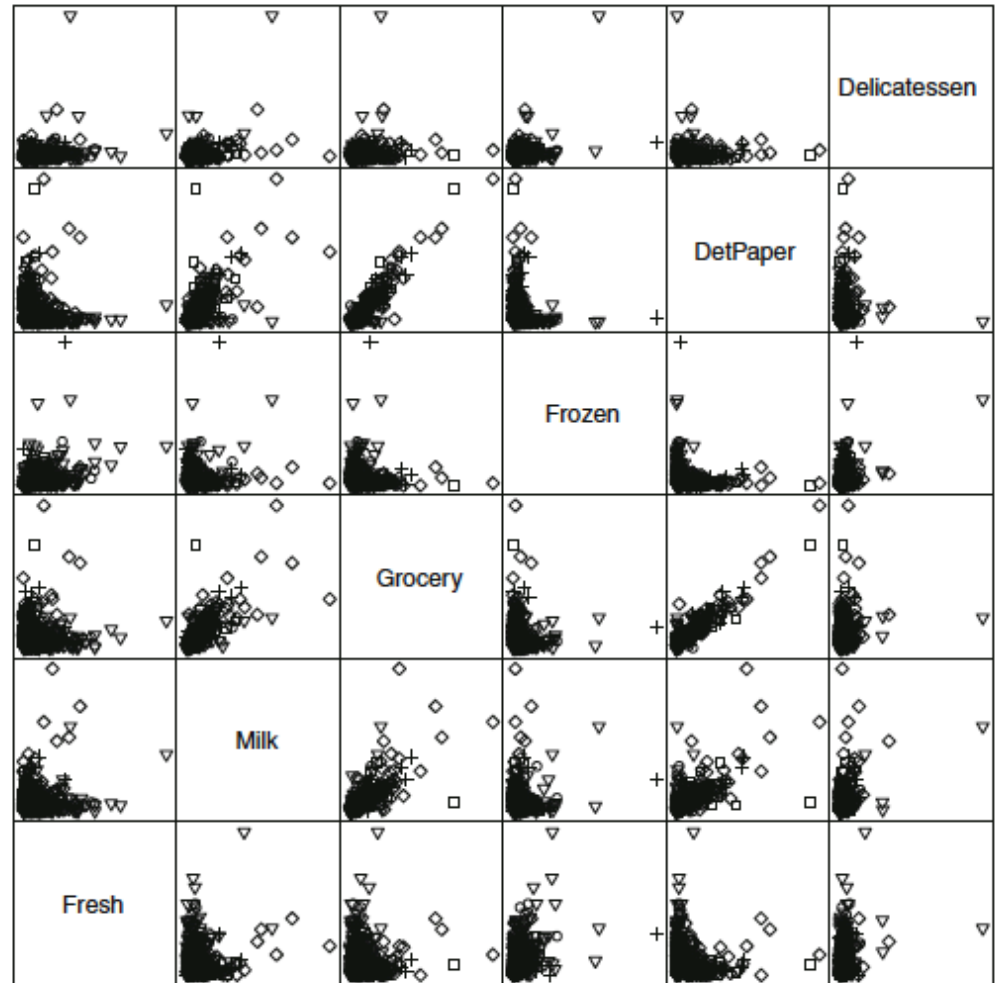✳ Region (Region 1, 2 &3) (Lisbon 77, Oporto 47, Other 316)

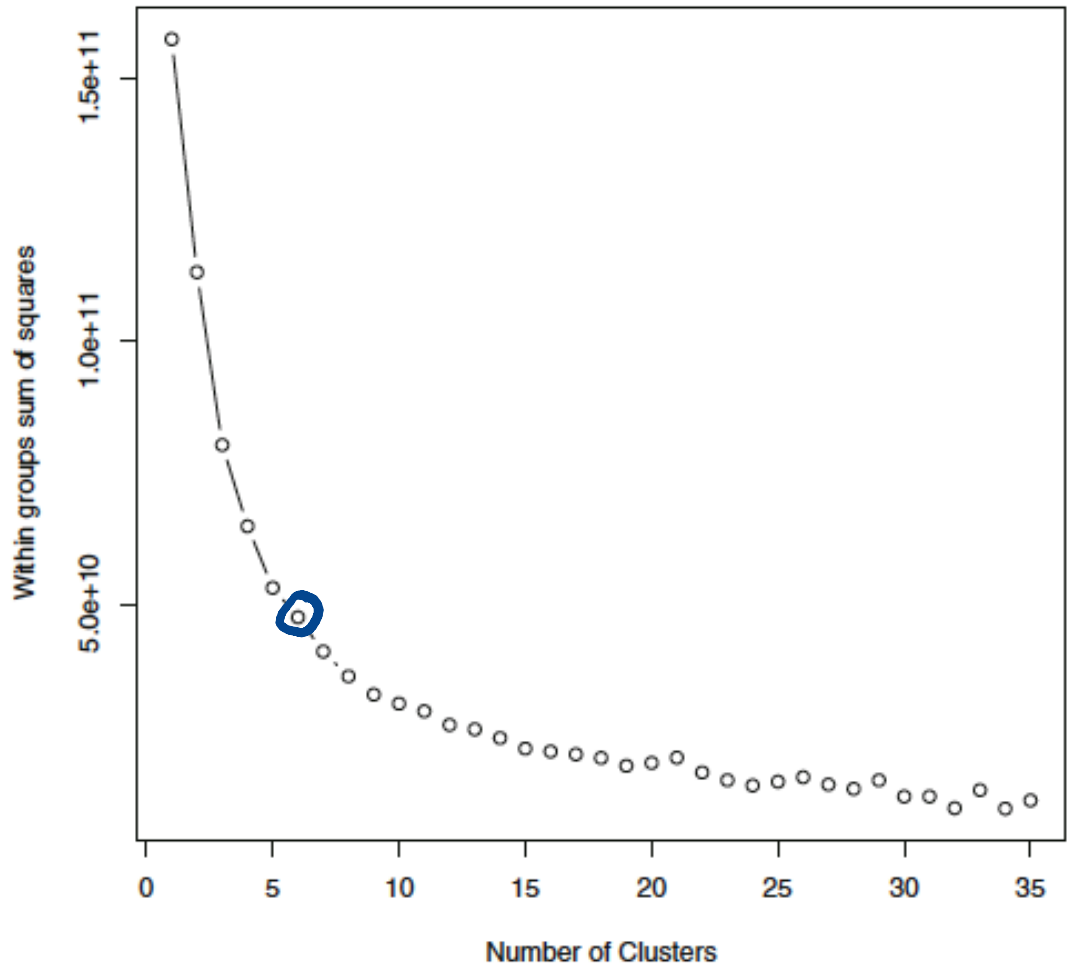# Lisbon, Portugal

# 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
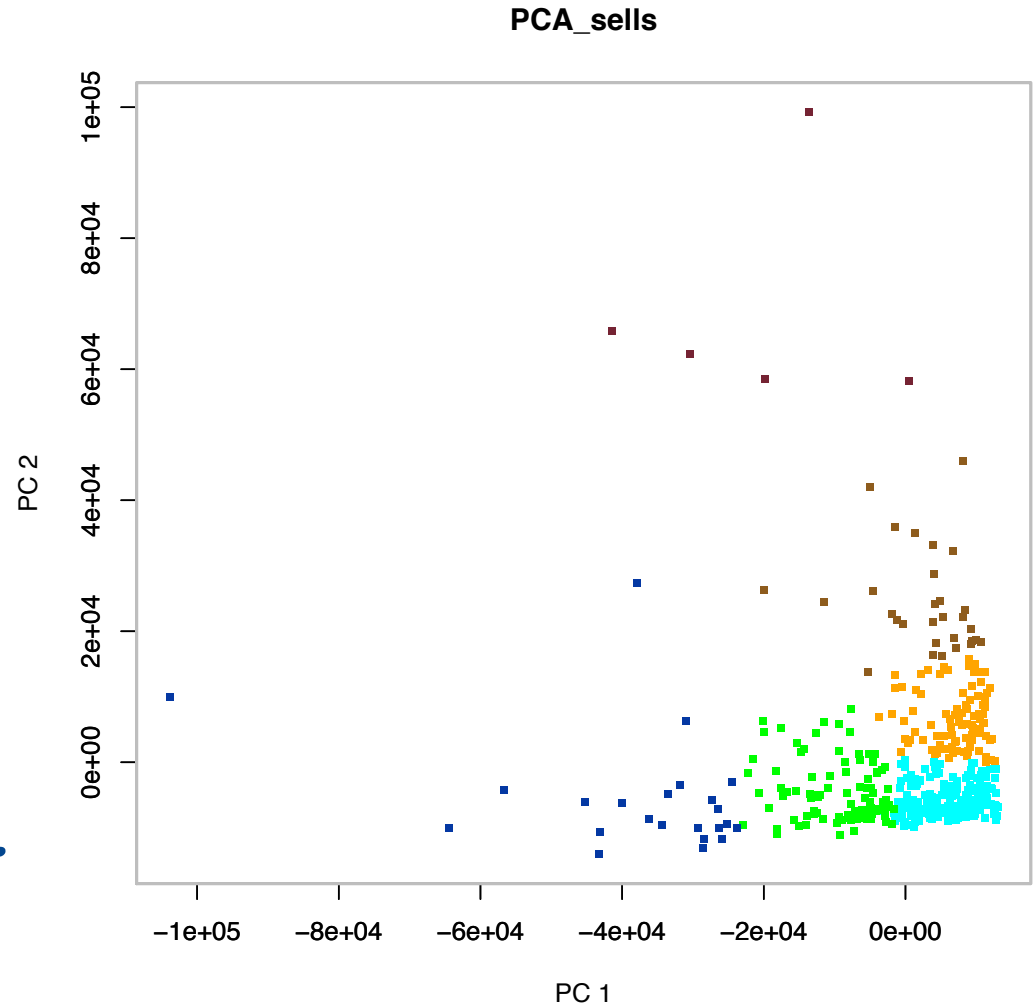
# Do kmeans and choose k through the cost function

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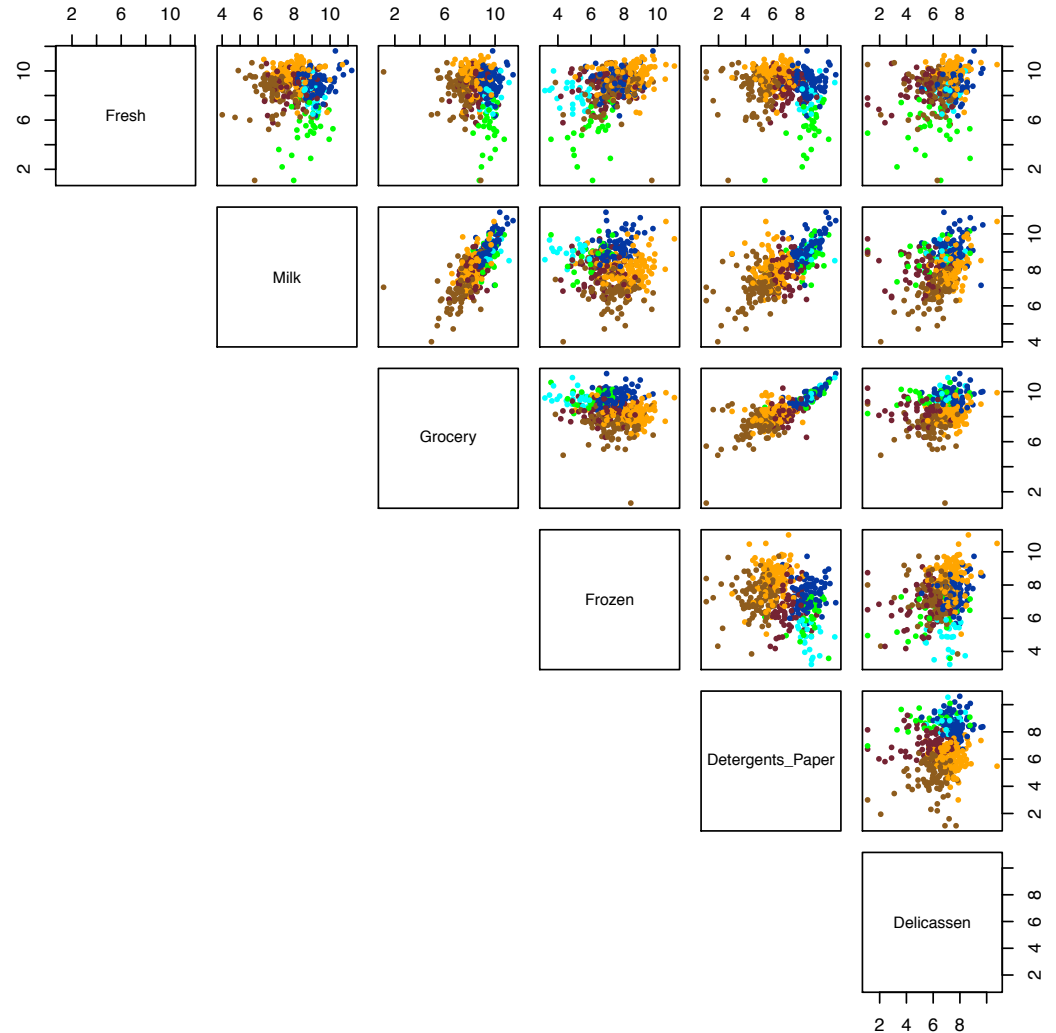
# Visualization of the data (PCA)

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

✳ Data points show large range of dynamics!
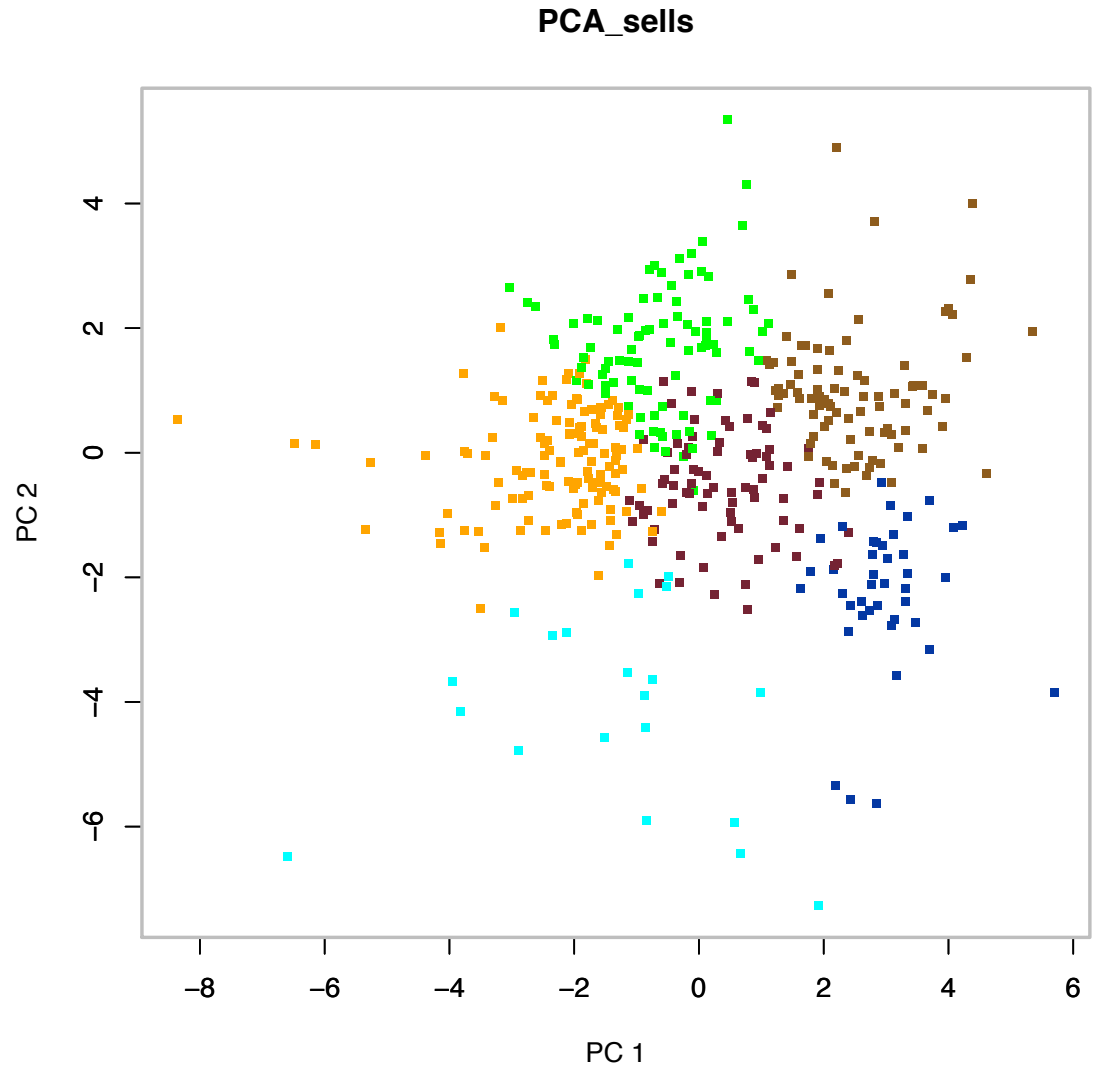
*each dot is one customer*



PCA_sells

# Do log transform of the data

- Log transform the data

- Do scatter plot matrix after the log transform

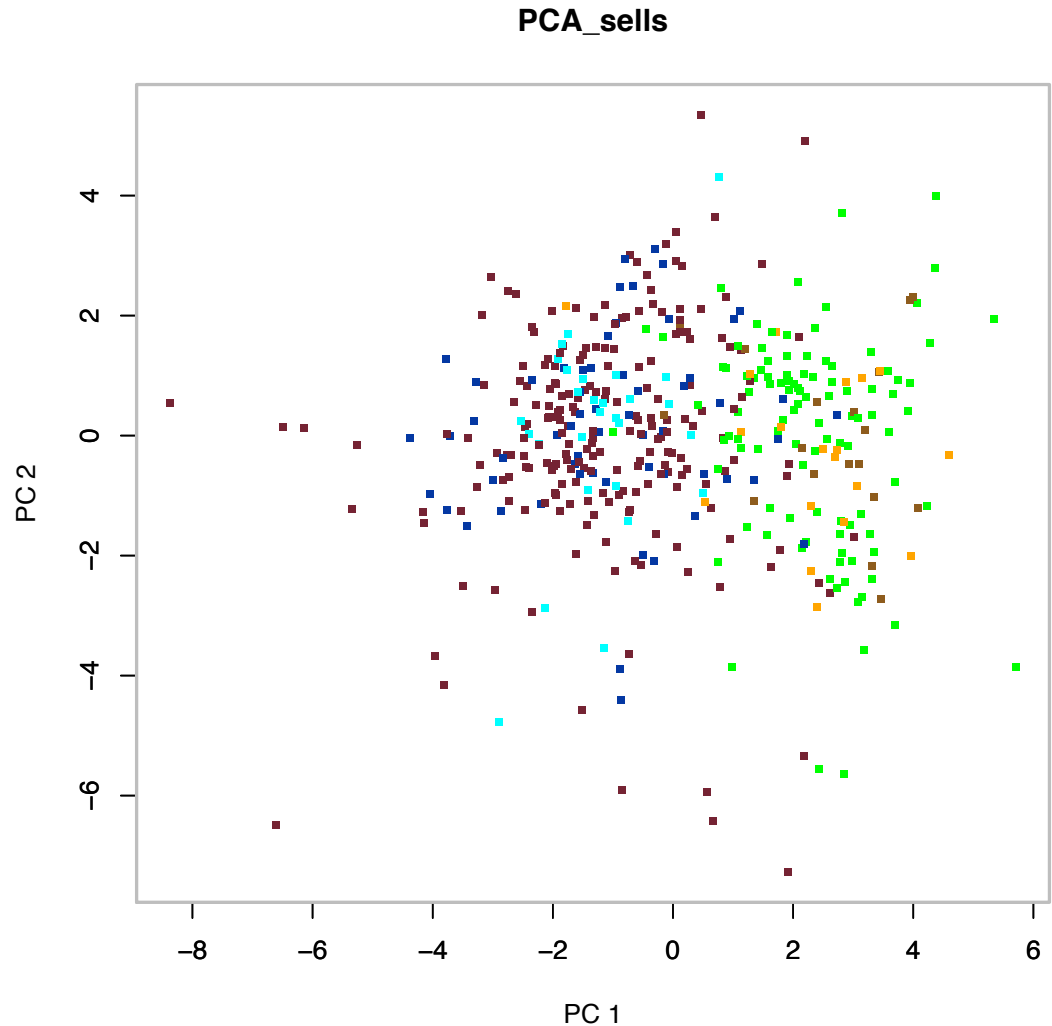- Do the kmeans and color the clusters identified by k-means

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



PCA_sells

# PCA after log transformation

Colors show the **Channel-region labels**
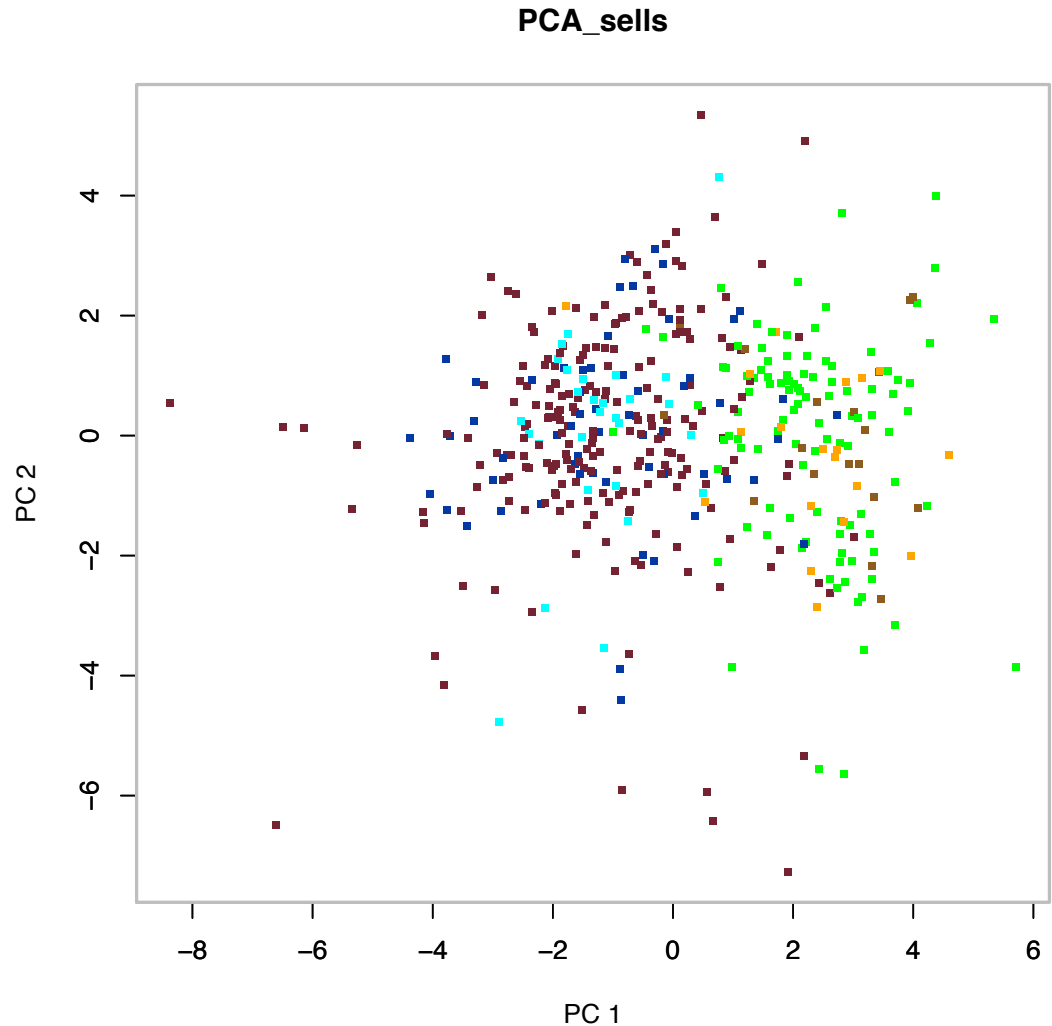
What does this tell us?



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

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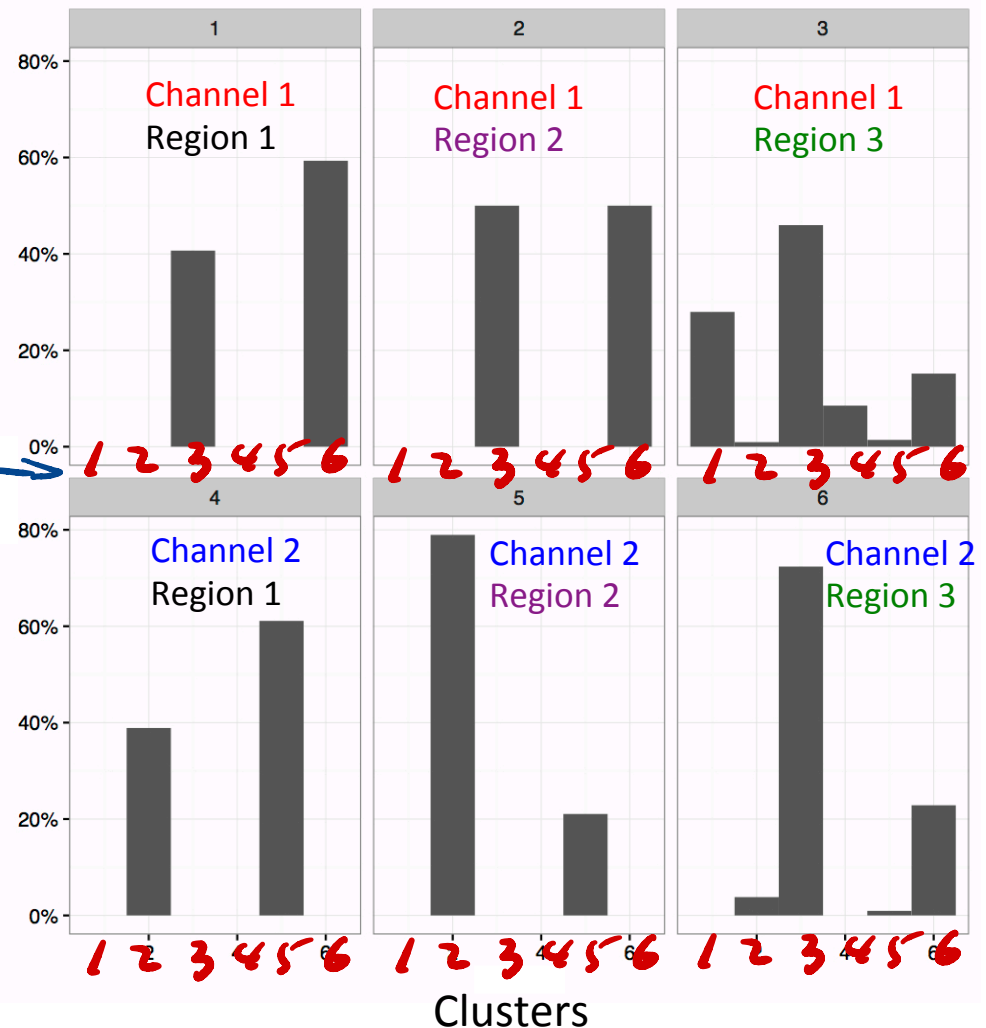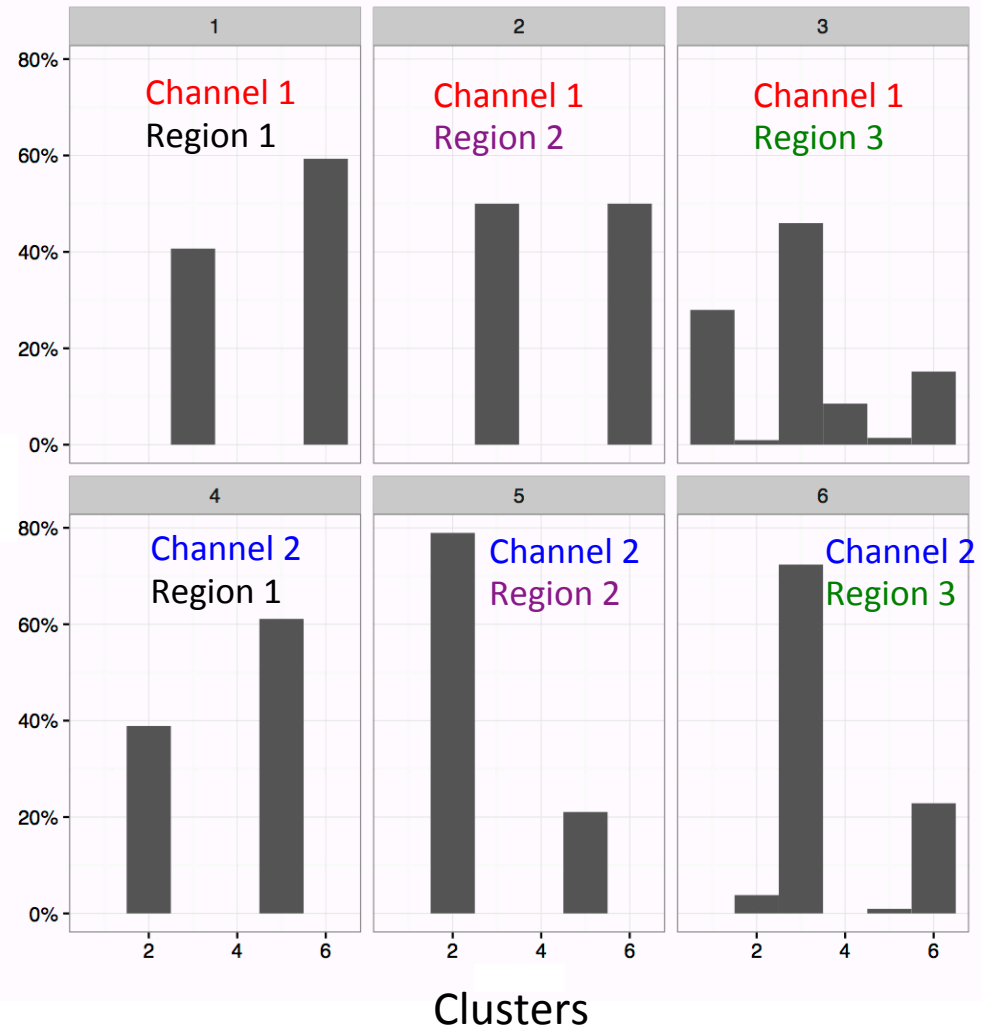
✳ **Is it enough to plot the percentage?**

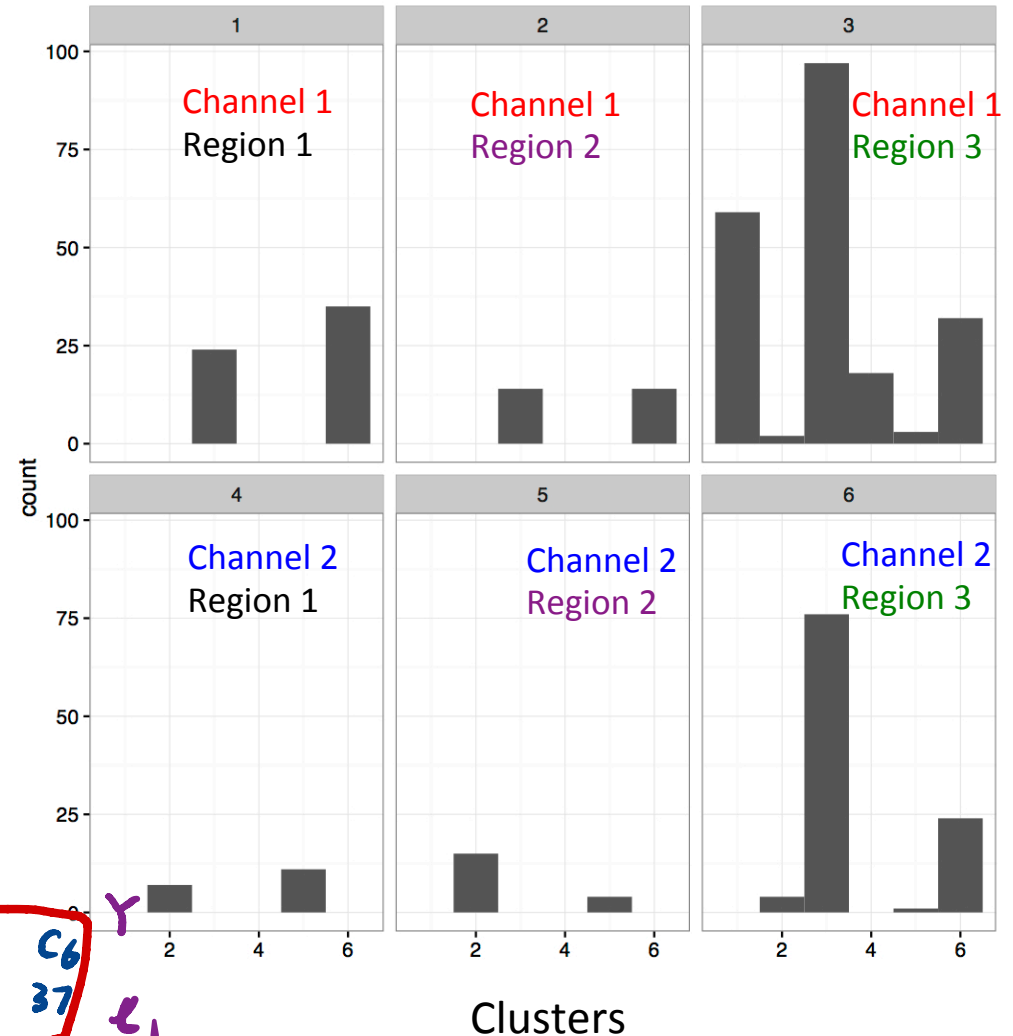# Cluster center histogram of the Portugal grocery spending data

- 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**



Clusters

$C_1$  $C_2$  $C_3$  $C_4$  $C_5$  $C_6$
0      0      24     0      0      37

$Y$

$L_1$

# 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. 75%

D. 90%

# 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

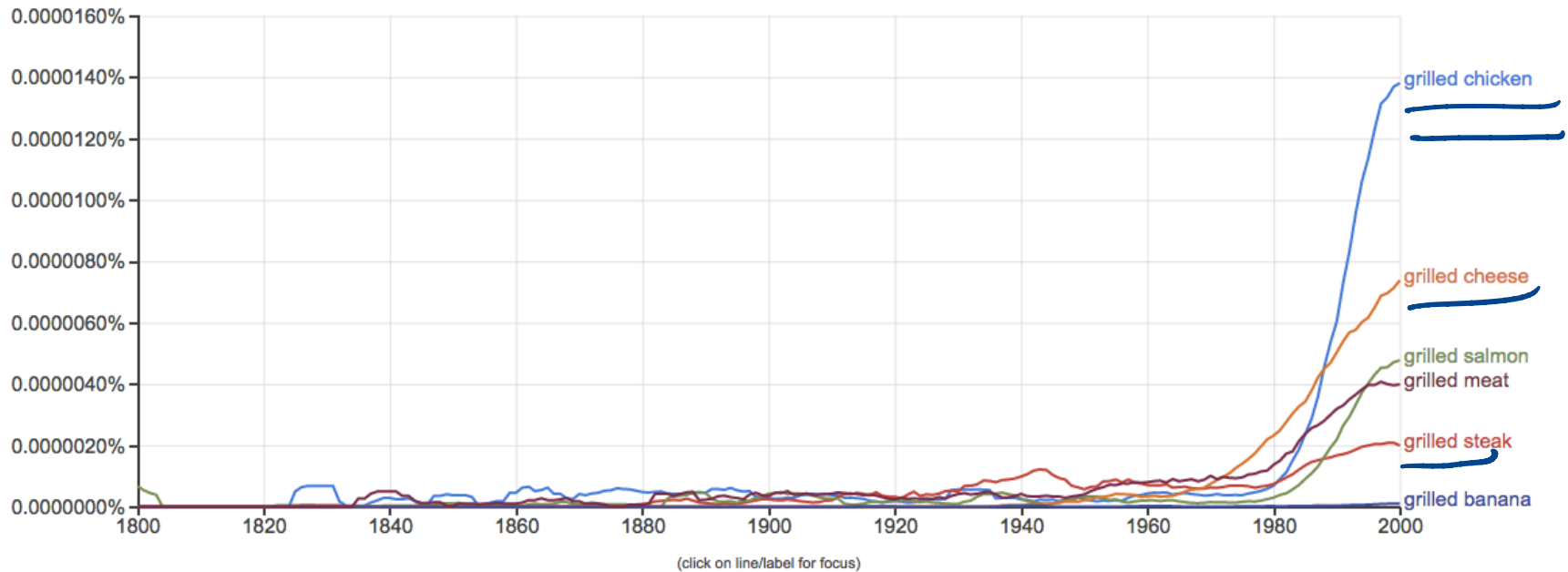I had a glass of wine with my grilled _____

*cheese*

*steak.*

*fish*

*shoe?*

*vege.*

*chicken*

# 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**.

✳ Such dependence is called **chain dependence**

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 $X_{t-1}$

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

Markov property

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

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

= C

# Coin example

✳ Toss a fair coin until you see two <u>heads</u> in a row and then stop, what is the probability of stopping after exactly **n** flips?
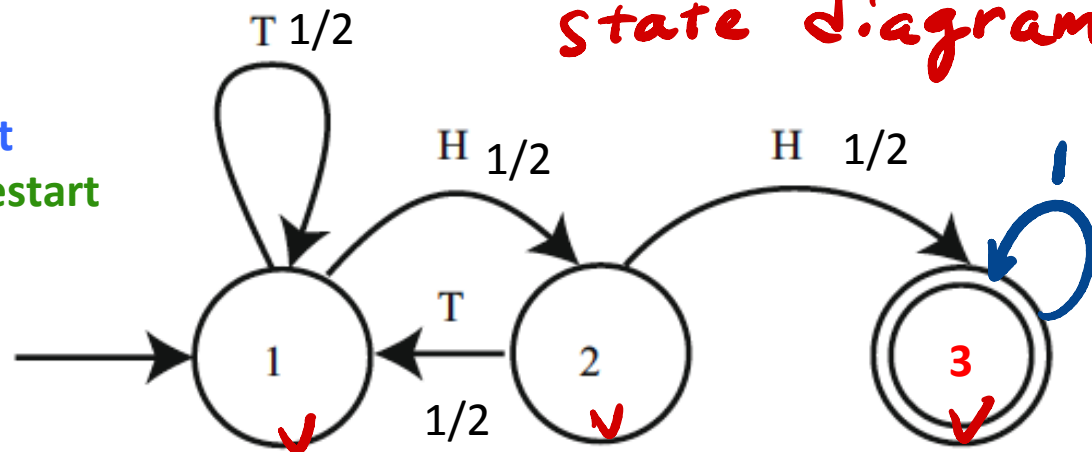
✱ ✱✱ HH

n

$$P(n = n_o) = ?$$

Geometric
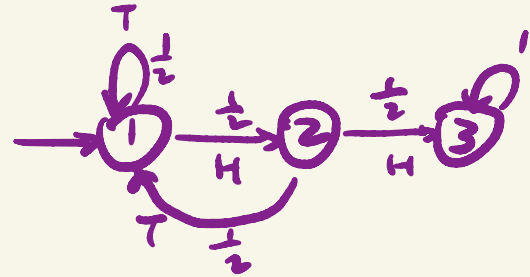
T T T... H

state diagram

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

T 1/2

H 1/2     H 1/2

T

1/2

1        2        3

$N =$ #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:



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

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

This part can be any !!

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

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

HHH

$$* H H$$
$$T H H$$

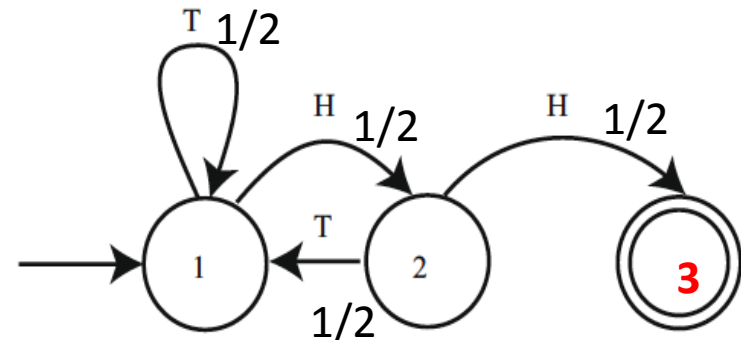$$* * H H$$
$$\{ T T H H$$
$$H T H H$$

$$P(n = n_0) =$$

$n\uparrow$

# The model helps form recurrence formula

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

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

$n$

※ 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

$T \quad \times \quad \times \quad \times \quad \times$

$n-1$

$n$

$HT \quad \times \quad \times \quad \times \times$

$n-2$

$P_{n-2} : P(\text{finish in } n-2 | HT)$

※ So we can derive a recurrence relation

$$p_n = \frac{1}{2}p_{n-1} + \frac{1}{4}p_{n-2}$$
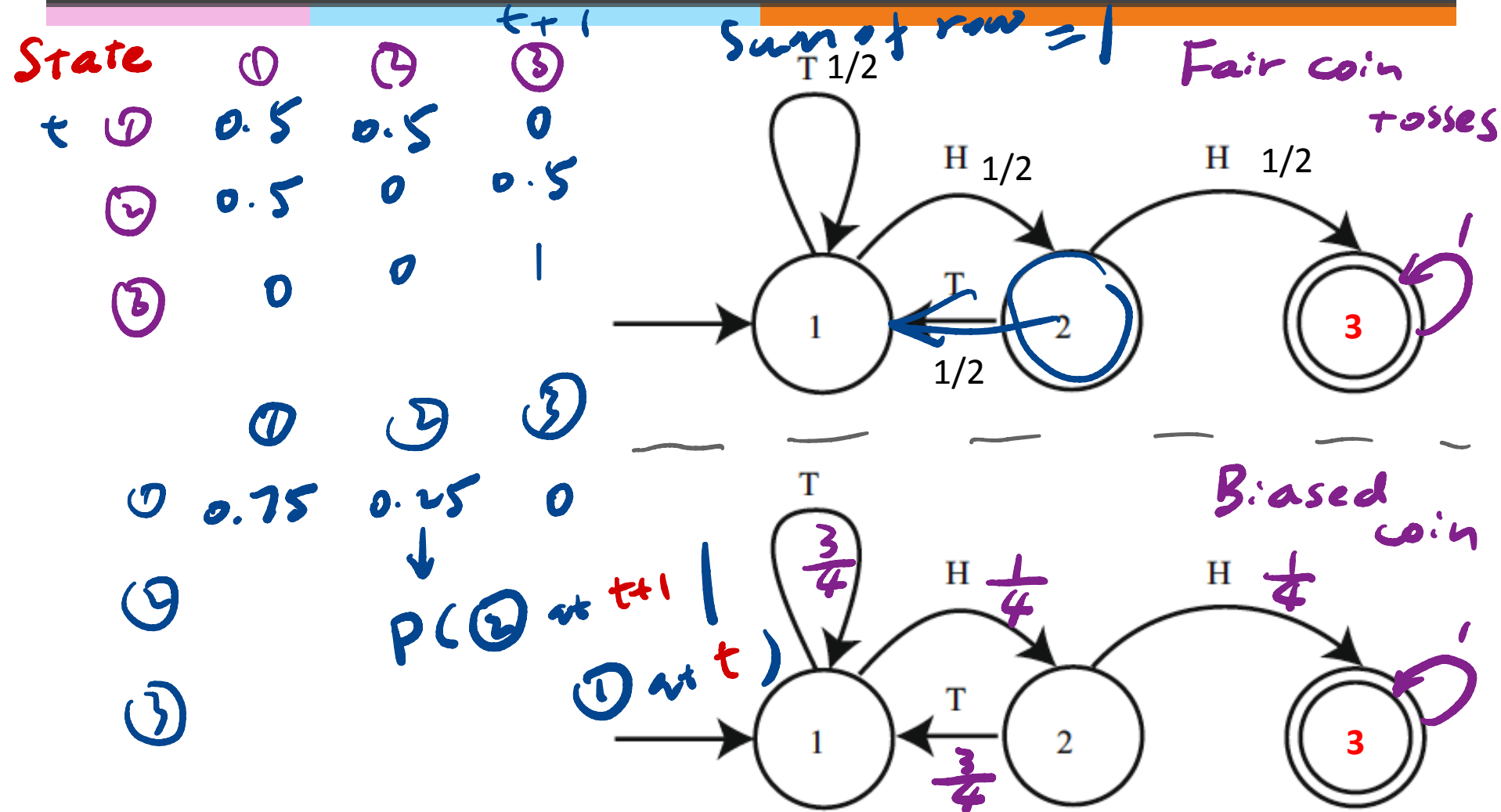
$P_{n-1} : P(\text{finish with } n-1 | T)$

P(T)       P(HT)

T 1/2

H 1/2      H 1/2

T

1      2      **3**

1/2

State

$t+1$

Sum of row = 1

Fair coin tosses

| State | ① | ② | ③ |
|---|---|---|---|
| t ① | 0.5 | 0.5 | 0 |
| ② | 0.5 | 0 | 0.5 |
| ③ | 0 | 0 | 1 |



T 1/2

H 1/2

H 1/2

T

1/2

| | ① | ② | ③ |
|---|---|---|---|
| ① | 0.75 | 0.25 | 0 |
| ② | | | |
| ③ | | | |

$P(②$ at $t+1$ | ① at $t$)

Biased coin

T $\frac{3}{4}$

H $\frac{1}{4}$

H $\frac{1}{4}$

T $\frac{3}{4}$

# 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.
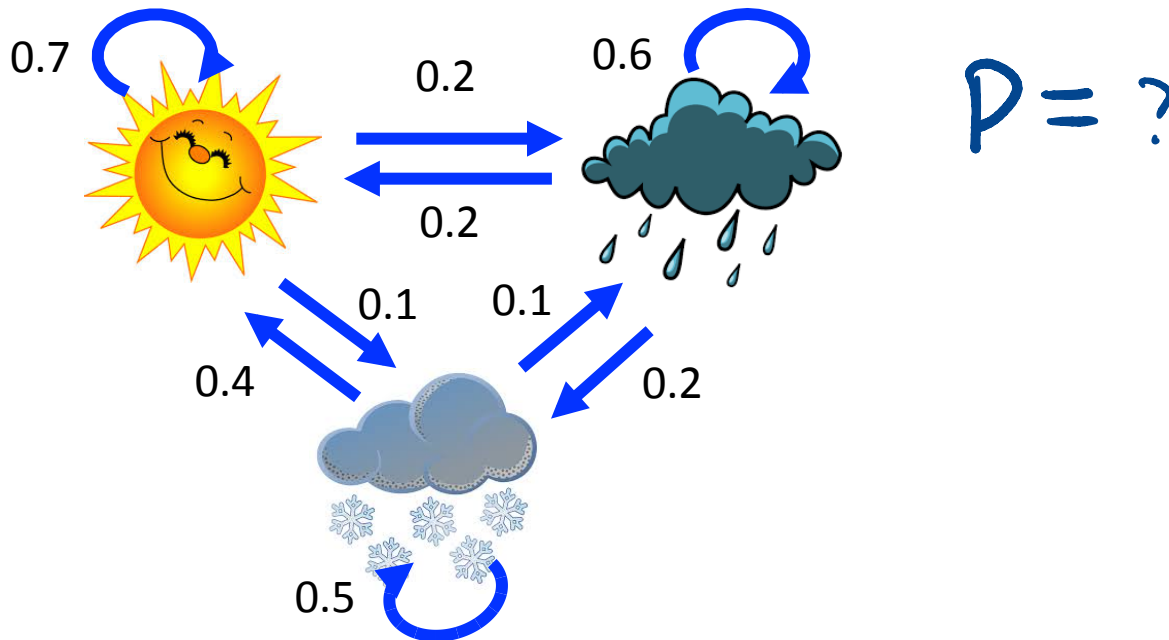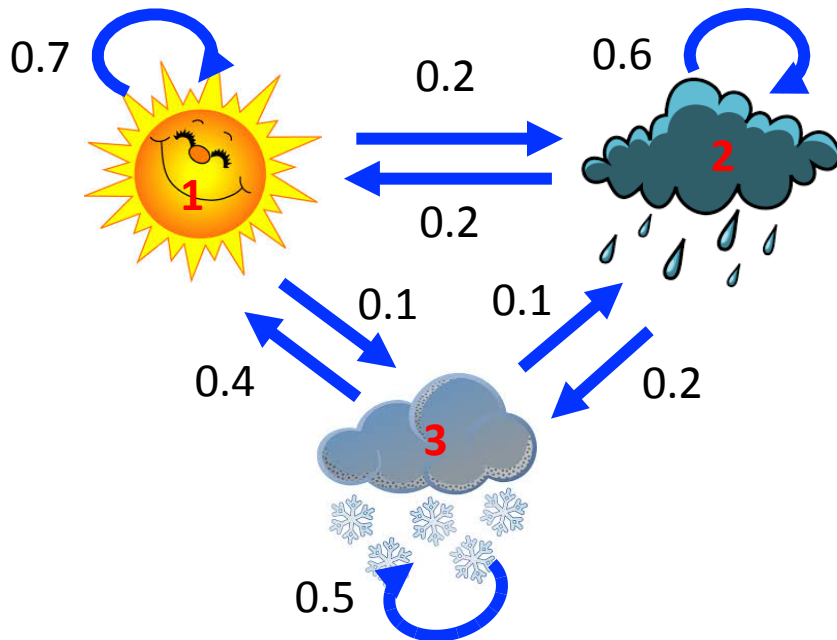


$$P = ?$$

# 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

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

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

$$\pi_0 = \begin{matrix} \text{Sunny} & \text{Rainy} & \text{Snowy} \\ [ \quad 0 & 1 & 0 \quad ] \end{matrix} \qquad \tau$$

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

P ( Snowy ) for next day ?        $t+1$

$$\pi_1 = \pi_0 P$$

$$= [ \; 0 \; 1 \; 0 \; ] \begin{bmatrix} \quad \end{bmatrix} = [ \begin{matrix} \text{Sunny} & \text{Rainy} & \text{Snowy} \\ & & \uparrow \\ & & ? \; ] \\ & & \uparrow \\ & & 0.2 \end{matrix}$$

$$\pi_1 = \pi_0 P$$

$$\pi_2 = \pi_1 P \quad ;$$

$$\pi_i = \pi_0 P^i$$

# 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 next time

*See You!*