

# Cadenas de Markov + Geometría de la Probabilidad + Estadística

ITAM

Clase 18 Curso Propedéutico  
2016/06/27

# Cadenas de Markov

- El estado próximo depende sólo del estado actual.
- Más precisamente, las *probabilidades de transición* están determinadas por el estado actual.
- **Motivación inicial:** un modelo para estudiar el lenguaje.



A. A. Markov (1886).

- Para conocer las probabilidades del futuro solo necesitas conocer el presente!

*Matemáticamente...*

- Si  $S$  son todos los posibles “estados” de la cadena (los valores que pueden tomar las  $X_t$ ) la propiedad de Markov dice:

$$P(X_t = x | X_0 = x_0, \dots, X_{t-1} = x_{t-1}) = P(X_t = x | X_{t-1} = x_{t-1})$$

Cadenas homogéneas

$$P(X_t = j | X_{t-1} = i) = P(X_1 = j | X_0 = i)$$

Más general

$$P(X_k = j | X_s = i) = P(X_{t+k} = j | X_t = i)$$

No ejemplo :

# Qué es una cadena de Markov

- Es tres cosas:

1. Un conjunto de estados  $S$
2. Una distribución inicial para  $X_0$
3. Un “Kernel” de transición. Unas probabilidades de transición por cada estado actual.

$S$  finito.... El kernel puede representarse como matriz  
 $P = (P_{ij})$

$$P_{ij} = P(X_{t+1}=j | X_t=i)$$

# Distribución de $X_n$ dado $X_0$

- La vez pasado vimos que  $P^n$  representa las probabilidades de transición en  $n$  pasos.
- Si  $\pi$  es un vector tal que  $\pi_j = P(X_0 = j)$  entonces tenemos que

$$P(X_1 = j) = \sum_i \pi_i P_{ij}$$

Por lo tanto  $\pi^\top P$  es un vector que representa la distribución de  $X_1$ .

En general,  $\pi^\top P^n$  es la distribución de  $X_n$

# Teorema ergódico

- El resultado más fundamental de las cadenas de Markov establece que si la cadena es “**regular**” entonces en el largo plazo no importa el estado inicial, la distribución de  $X_n$  no depende de  $X_0$ .
- Tal distribución se conoce como **distribución estacionaria/límite** y cumple.

$$\hookrightarrow \pi = P^T \pi \quad \hookrightarrow$$

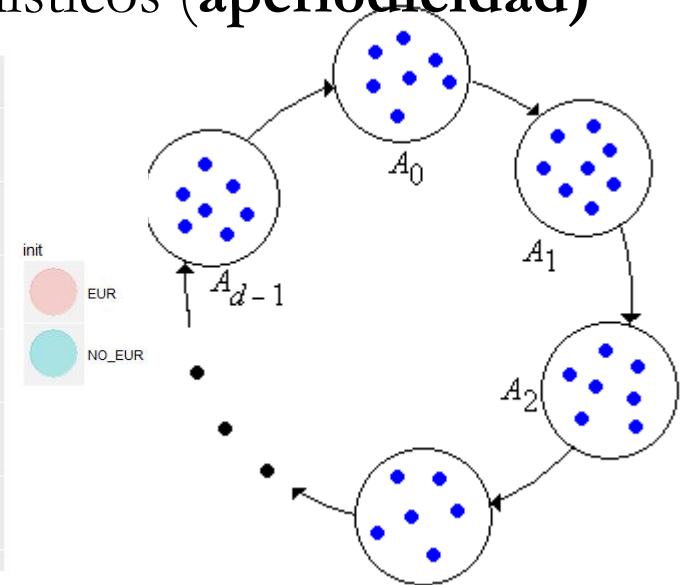
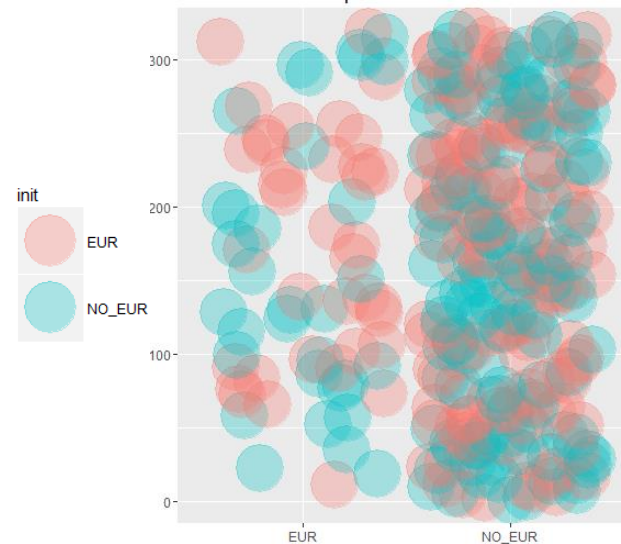
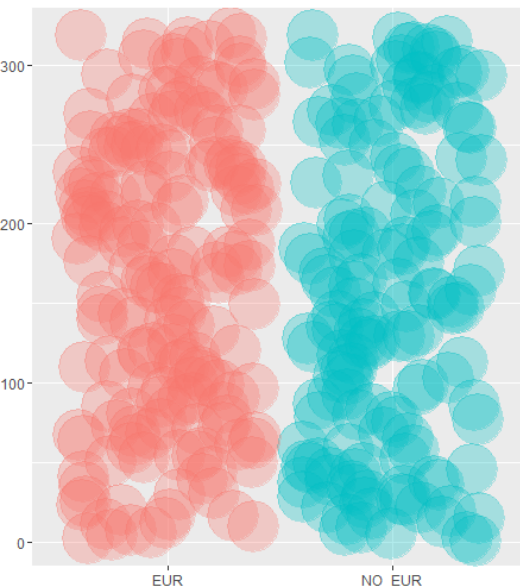
$$\pi = \lim_{k \rightarrow \infty} \pi^{(k)}$$

$$\pi_j^{(k)} = P(X_j = k)$$

~~Y como consecuencia:~~

- Aunque no podemos verlo con detalle, hay que recordar un caso particular en el que el *teorema ergódico* se cumple siempre que:

1. Todos los estados están conectados en el sentido de que existe alguno camino de estados tal que siempre se puede llegar de un estado a otro y de regreso (**irreducibilidad**).
2. El tiempo esperado de retorno a cada estado es finito (**recurrencia positiva**)
3. La cadena no tiene *ciclos* determinísticos (**aperiodicidad**)



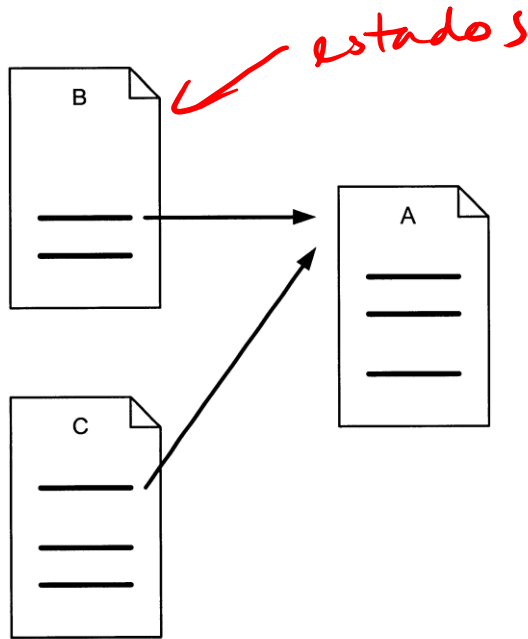


# Google PageRank



- **Motivación:** El problema de un “buscador”: un usuario introduce una *búsqueda*, las páginas web se encuentran *indexadas*.
- $S$  representa el conjunto de páginas web que responden a la búsqueda
- **Objetivo:** En qué orden listar  $S$  al usuario.

- Necesitamos un método de **rankeo**!



**FIG. 1**

**Approach:** Define una cadena de Markov que represente cómo se mueve un “usuario” a través de páginas web y ver que estados son los más visitados en el largo plazo.

# ¿Cómo construir $P$ ?

- Una propuesta:

Estimaremos las probas de transición basados en contar qué porcentaje de los links dirigen a cada página.

clicks  
historial etc ← construir probas de transición

*Posibles Problemas par que no exista la distribución límite?????*

Podría no ser irreducible

- Los estados podrían estar desconectados!!!!
- Para eso se agrega un pequeño “ruido” que corresponde a moverse completamente aleatoriamente entre cada página web.
- Se escoge un número  $\alpha \in (0,1)$  (e.g  $\alpha \approx .85$ ) y se pone

$$p_{ij} = \alpha \frac{l_{ij}}{\sum_k l_{ik}} + (1 - \alpha) \frac{1}{n}$$

Donde  $n$  es el # total de páginas web y  $l_{ij}$  el total de link en la página  $i$  con destino  $j$ .

# **LA GEOMETRÍA (Y ÁLGEBRA LINEAL) DE LA PROBABILIDAD**

## Recordatorios y Observaciones

1. Espacio vectorial: conjunto cuyos elementos sabemos sumar y multiplicar por un escalar, y es cerrado bajo estas operaciones

Obs 1: Sabemos sumar variables aleatorias

$$X, Y: \Omega \rightarrow \mathbb{R}$$

$$(X+Y)(\omega) = X(\omega) + Y(\omega)$$

También se pueden multiplicar por escalar

$$X: \Omega \rightarrow \mathbb{R} \quad a \in \mathbb{R}$$

$$(aX)(\omega) = a(X(\omega))$$

(Esto es cierto para funciones, no solo v.a.s)

Las variables aleatorias son un espacio vectorial

(de dimensión infinita)  
no tiene base

2. Propiedades del producto punto

operador  
bilineal

$$\left\{ \begin{array}{l} \bullet) \langle X, X \rangle \geq 0 \quad \langle X, X \rangle = 0 \Leftrightarrow X = 0 \\ \bullet) \langle X, Y \rangle = \langle Y, X \rangle \\ \bullet) \langle X + Y, Z \rangle = \langle X, Z \rangle + \langle Y, Z \rangle \\ \quad \vee \langle \alpha X, Y \rangle = \alpha \langle X, Y \rangle \end{array} \right.$$

en  $\mathbb{R}^n$

$$\langle X, Y \rangle = X^T Y = X \cdot Y = \sum_{i=1}^n X_i Y_i$$

Obs 2: En variables aleatorias  $\langle X, Y \rangle = \mathbb{E}(XY)$  cumple las mismas

propiedades salvo  $\langle X, X \rangle = \mathbb{E}(X^2) = 0 \Leftrightarrow X = 0$ . En probabilidad  
 $\mathbb{E}(X^2 = 0) \Rightarrow P(X=0) = 1$  pero podría existir  $\omega$  con  $X(\omega) \neq 0$  pero  
 $P(X \neq 0) = 0$ , entonces es 'casi' como  $X=0$ . Formalmente  
 se dice  $X=0$  casi seguramente. Por lo tanto,  $\langle X, Y \rangle = \mathbb{E}(XY)$   
 es un producto punto.

$$\cos(\theta(\vec{x}, \vec{y})) = \frac{\langle X, Y \rangle}{\sqrt{\langle X, X \rangle \langle Y, Y \rangle}} \Rightarrow \text{En v.a.s} \quad \theta(X, Y) = \cos^{-1} \left( \frac{\mathbb{E}(XY)}{\sqrt{\mathbb{E}(X^2) \mathbb{E}(Y^2)}} \right)$$

¿Existen ángulos entre variables aleatorias!

(Las v.a.s son  
espacios de  
Hilbert)

3. Relación entre norma y producto punto

$$\langle X, X \rangle = \|X\|^2$$

En v.a.s  $\|X\| = \sqrt{\mathbb{E}(X^2)} = \sqrt{\langle X, X \rangle}$

En  $\mathbb{R}^n$

$$\|X\|_2 = \sqrt{\sum X_i^2} = \sqrt{\langle X, X \rangle}$$

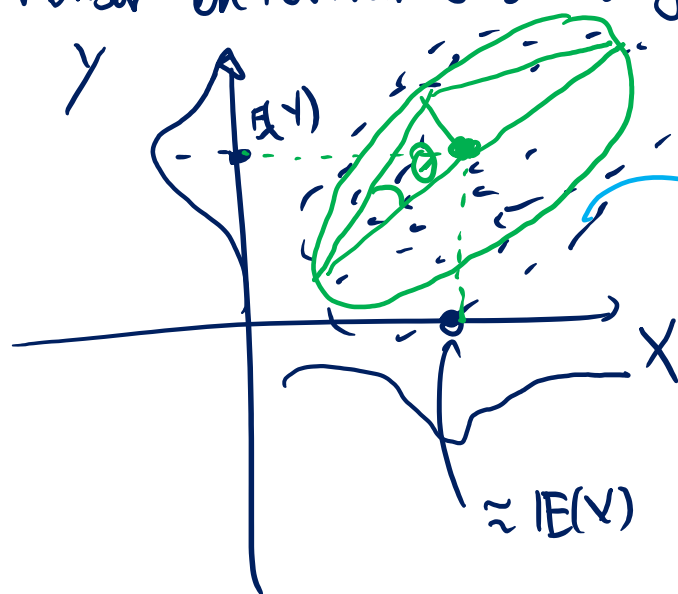
4. Este approach nos da una nueva versión de la covarianza y la varianza

$$\text{Cov}(X, Y) = \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y))) = \langle X - \mathbb{E}(X), Y - \mathbb{E}(Y) \rangle$$

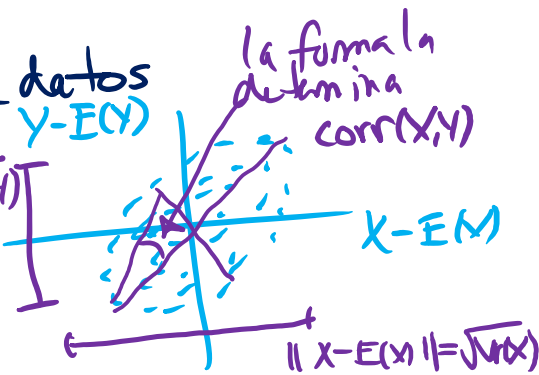
$$\text{Var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \|X - \mathbb{E}(X)\|^2$$

Podemos entonces pensar las covarianzas y varianzas en términos de "ángulos" y "normas"

Pensar en términos de diagramas de dispersión de datos



← cada punto es la coordenada  $(x_i, y_i)$  de un individuo



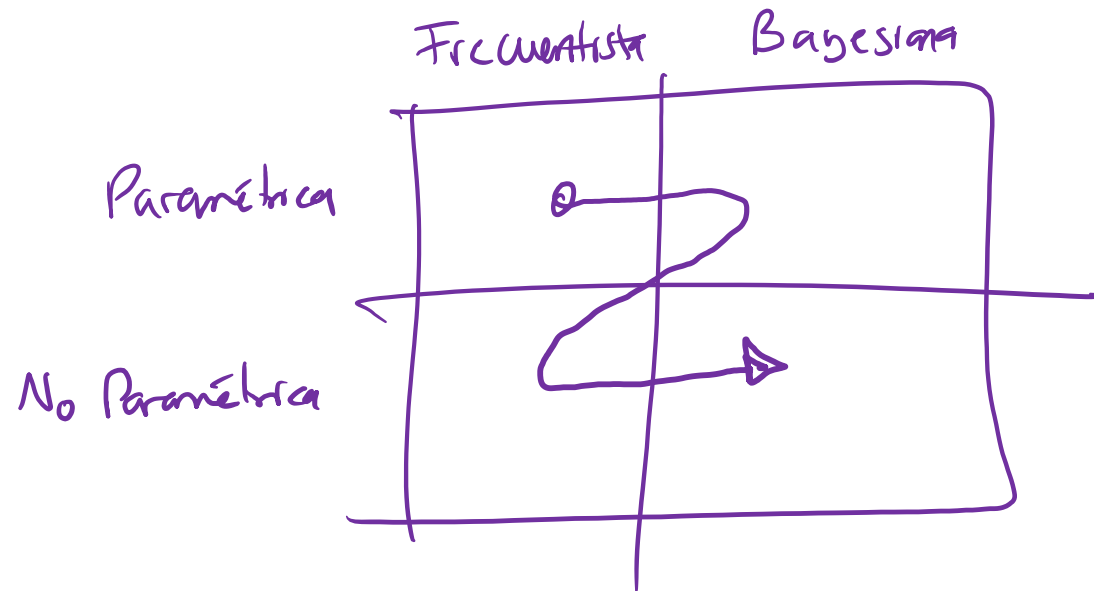
La elongación del diagrama depende de  $\text{Cov}(X, Y)$

$$\begin{aligned} \text{Corr}(X, Y) &= \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\langle X - \mathbb{E}(X), Y - \mathbb{E}(Y) \rangle}{\|X - \mathbb{E}(X)\| \|Y - \mathbb{E}(Y)\|} \\ &= \cos(\angle(X - \mathbb{E}(X), Y - \mathbb{E}(Y))) \end{aligned}$$

Continuara....



# ESTADÍSTICA

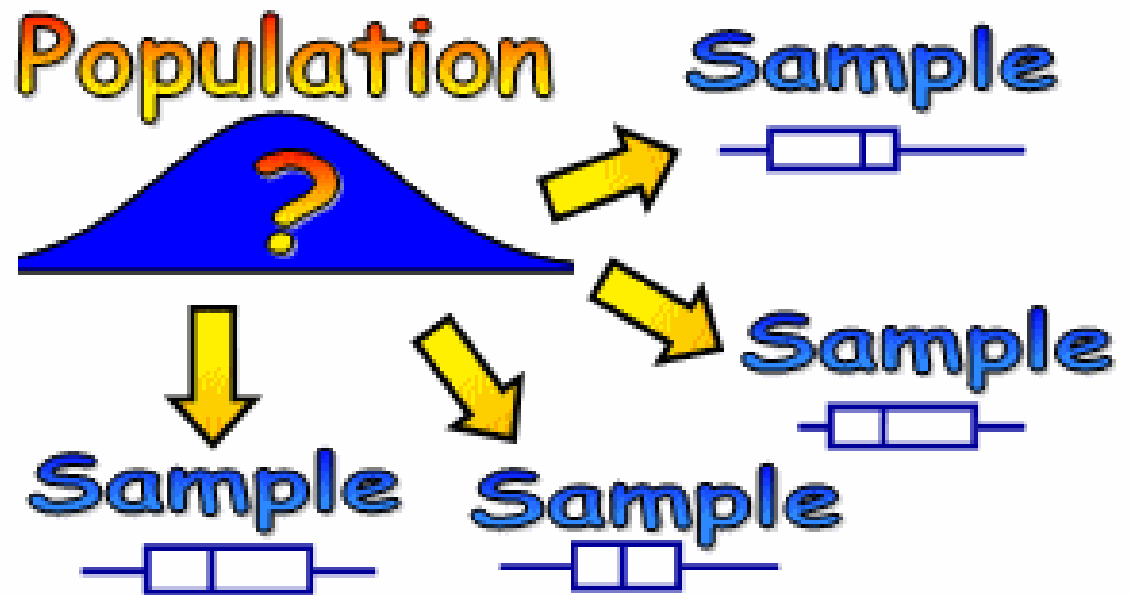


# Muestras

Es un supuesto que la estadística básica cambiar por "intercambialidad" pero tienen permiso de hacerlo.

- Una **muestra** es una colección de variables aleatorias  $X_1, \dots, X_m$  usualmente independientes con una misma distribución  $F_X$
- Usualmente en proba decimos que es una sucesión de v.a.i.i.d.
- Usar la palabra muestra o v.a.i.i.d. es usualmente un problema de *enfoque*.

- Vaidis vs Realizaciones
- $X_i$  vs  $X_i(\omega) = \underline{x_i}$



# Estadística Estimadores

- Una **estadística** es cualquier cantidad de interés acerca de una distribución  $F_X$  (*in fancy words: una funcional*)
- Dada una muestra  $X_1, \dots, X_n$ , un **estimador** es una función  $h(X_1, \dots, X_n)$ .
- Un estimador es una función que depende de los valores de la muestra y que su propósito es tratar de conocer el valor de alguna estadística.

# Ejemplos....

- *Teaser*: ¿Cómo estimarían la esperanza de una distribución dada una muestra?

Queremos  $\int_{\mathcal{X}} x dF_X(x) = \int_{\mathcal{X}} x f_X(x) dx = \mathbb{E}(X_i) \triangleq \eta(F_X)$

de una muestra  $X_1, \dots, X_n$

$\hat{\eta}(X_1, \dots, X_n) = \frac{1}{n} \sum X_i$

función de la distrib

La estadística es la media

$\hat{\eta}$  estima  $\eta$

estimator estadística

¿cómo saber si un estimator es bueno/malo?

# Estadística Paramétrica

- Tenemos una muestra  $X = \{X_1, \dots, X_n\}$  con distribución común  $F_X^\theta$  y queremos aprender el “valor” de  $\theta$ . Queremos “decidir” que valores puede tomar  $\theta$ .
- Ej:
- $F_X^\theta \sim N(\mu, \sigma^2) \quad \theta = (\mu, \sigma^2)$
- $F_X^\theta \sim Ber(p) \quad \theta = (p)$

En la vida real nosotros acotamos a una familia paramétrica específica dependiendo del tipo de datos. La elección de familia paramétrica se llama el **modelo**.

# Inferencia Bayesiana vs Inferencia Frecuentista

Inferencia Frecuentista	Inferencia Bayesiana
Los parámetros $\theta$ que queremos conocer son tratados como NÚMEROS	Los parámetros $\theta$ que queremos conocer son tratados como VARIABLES ALEATORIAS para modelar nuestra incertidumbre.
NO puedo aportar conocimiento previo.	Puedo aportar conocimiento previo o “a priori” que influye en el resultado final.
Dada una muestra, me devuelve un NÚMERO que un estimador.	Dada una muestra me devuelve una nueva distribución para $\theta$ llamada distribución “a posteriori” que combina la “a priori” con los datos.
Máximo verosimilitud es un <i>ejemplo</i> de estimación frecuentista.	Hay un método de inferencia: <b>el método de Bayes</b> y combina la <i>verosimilitud</i> con la información a priori.

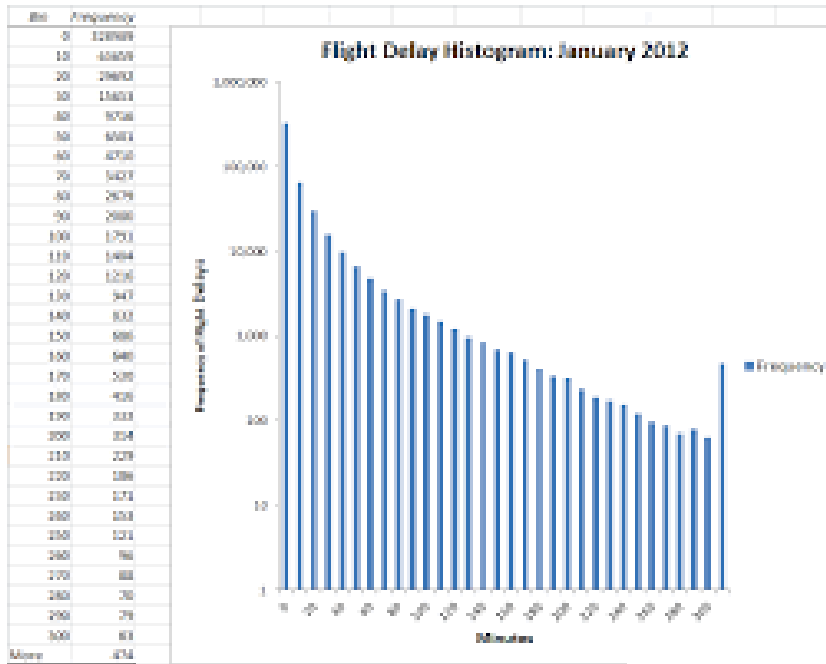
(*just to be clear...* El método de *máxima verosimilitud* es un método 100% frecuentista, pero la verosimilitud tiene un rol importante en la inferencia bayesiana...)

# Máxima verosimilitud

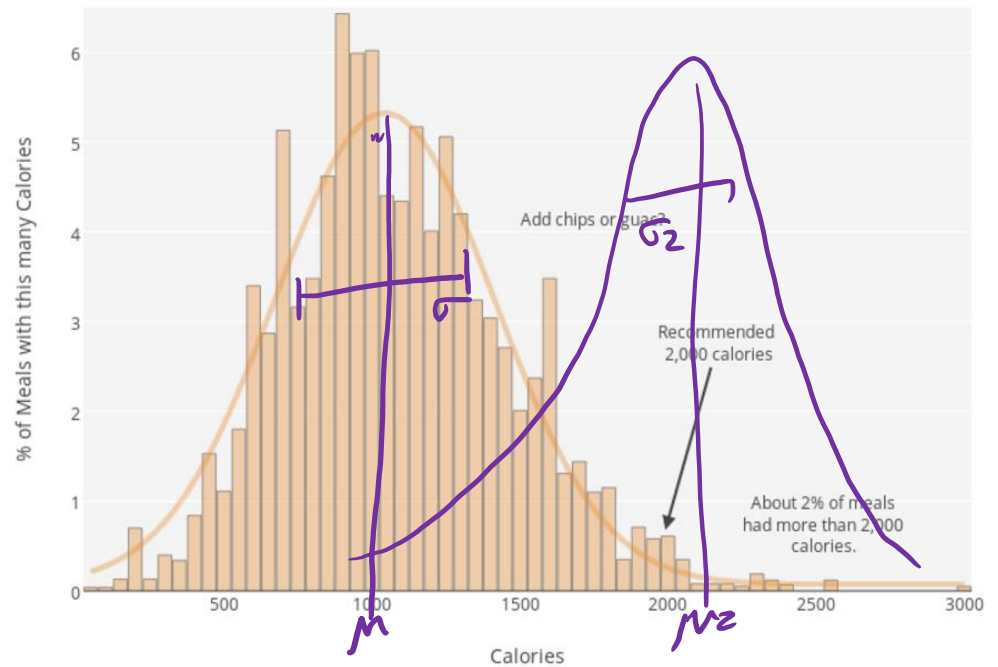
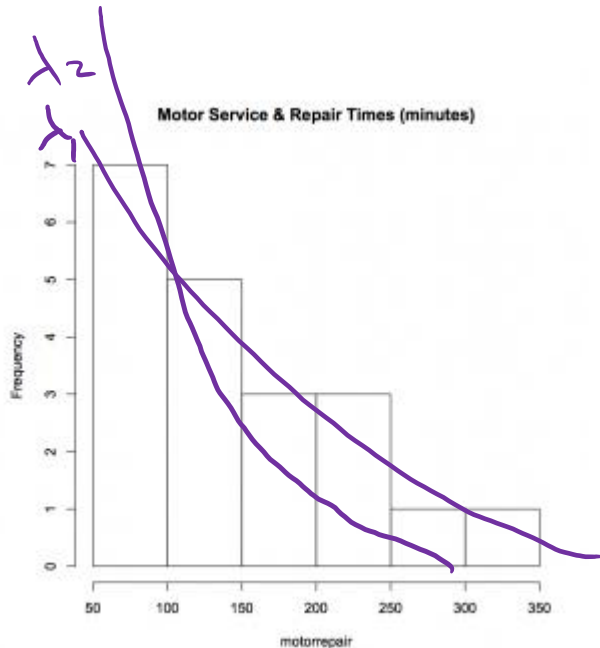
- Si  $F_X \sim \text{Ber}(p)$  y  $X_1, \dots, X_n$  es una muestra cuya realización fue:  
1,0,1,1,1,1,0,1,0,1,0,1,1,1,0,1,1,1,1
- ¿Qué  $p$  es más **creíble**?  $p = .3$  o  $p = .75$ ? Y maximizarla...

Nuestro objetivo es buscar una “medida” de qué tan creíble es cada  $p$  y buscar la  $p$  que maximice esta credibilidad: en esencia eso es máxima verosimilitud.



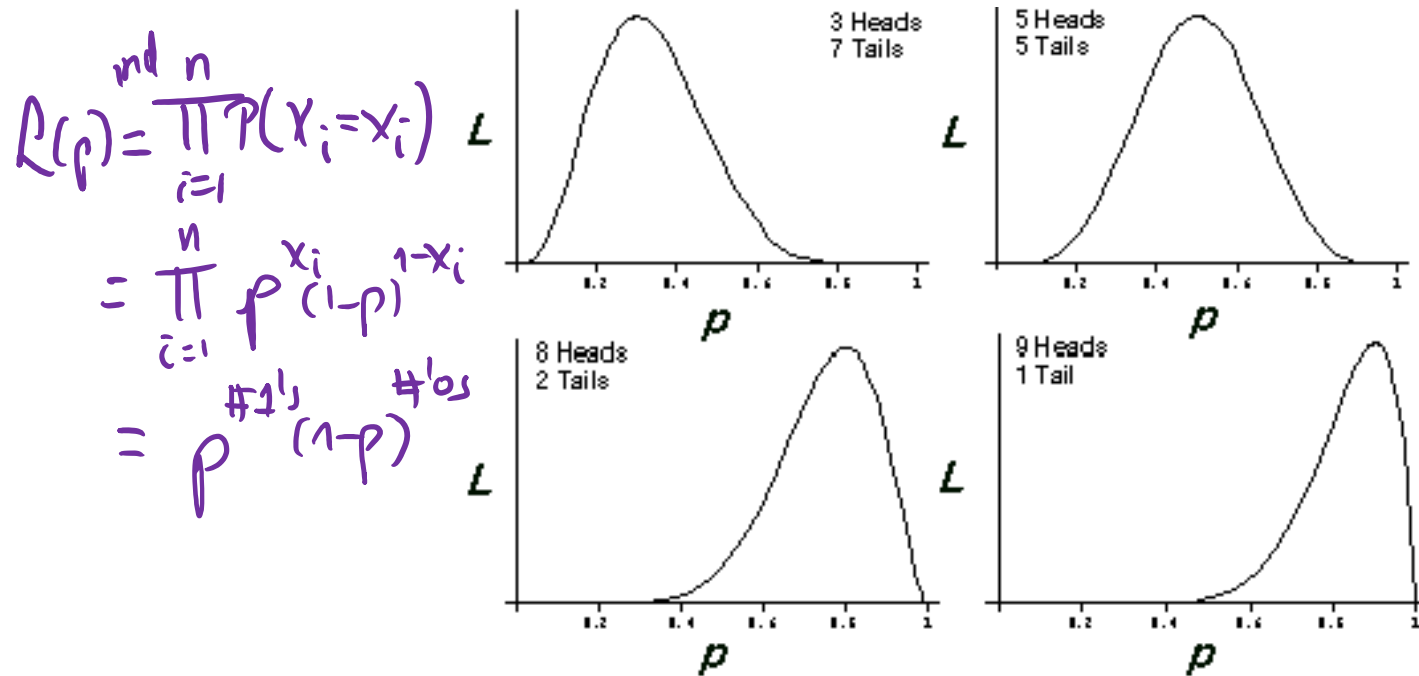


At Chipotle, How Many Calories Are You Consuming?



- Existe una elección de función de “credibilidad” que es muy natural:
- Si en el ejemplo anterior graficamos

$L(p) = P(X_1 = x_1, \dots, X_n = x_n; p)$  para distintos valores de  $p$  tenemos obtenemos algo que se ve así:



- La función de **verosimilitud** se define como:

$$L(\theta|X) = P(X_1 = x_1, \dots, X_n = x_n; \theta) \\ = \prod_i P(X_i = x_i; \theta)$$

Si la distribución es discreta solo hay que reemplazar la probabilidad por la función de densidad:

$$L(\theta|X) = f_X(X_1 = x_1, \dots, X_n = x_n; \theta) \\ = \prod_i f_X(X_i = x_i; \theta)$$

Método de Máxima verosimilitud

Inferencia clásica:

Dada una muestra  $X_1, \dots, X_n$  -  
 $\hat{\Theta}_{MLE}(X_1, \dots, X_n) = \underset{\Theta}{\operatorname{argmax}} L(\theta|X)$

# Inferencia Bayesiana

Teorema de Bayes:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)} \propto P(B|A)P(A)$$

Llevémoslo a variables aleatorias y pongan  $\theta$  en vez de  $A$  y su muestra  $X = \{X_1, \dots, X_n\}$  en vez de  $B$  y obtienen

$$f_{\theta|X}(\theta|X) = \frac{\overbrace{L(\theta|X)}^{P(X|\theta)} f(\theta)}{\int L(\tau|X) f(\tau) d\tau} \propto L(\theta|X) f(\theta)$$

en el caso continuo... pongan sumas y probabilidades en vez de integrales y densidades en el caso discreto y ya...

**“están mezclando la a priori con la verosimilitud”**

