# Distinct elements

## Anna Karlin
### With many slides by Luxi Wang, Shreya Jayaraman, Alex Tsun and Jeff Ullman

# Data mining

- In many data mining situations, the data is not known ahead of time.

- Examples:

  - Google queries

  - Twitter or Facebook status updates

  - Youtube video views

- In some ways, best to think of the data as an infinite stream that is non-stationary (distribution changes over time)
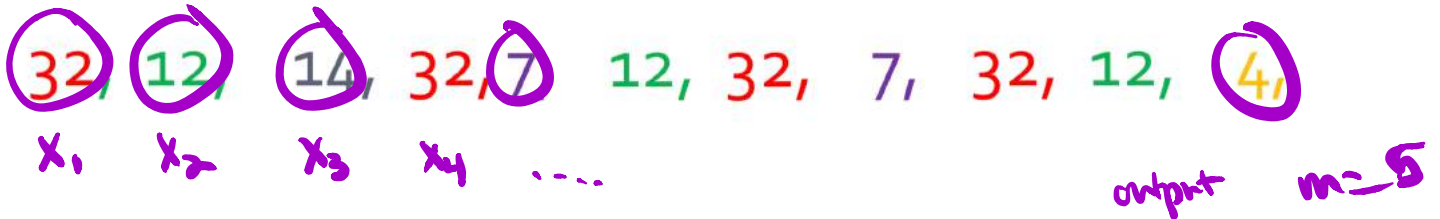
# Stream model

- Input elements (e.g. Google queries) enter/arrive one at a time.
- We cannot possibly store the stream.


Question: How do we make critical calculations about the data stream using a limited amount of memory?

# Problem

- Input: sequence of $N$ elements $x_1, x_2, \ldots, x_N$ from a known universe $U$ (e.g., 8-byte integers).

- Goal: perform a computation on the input, in a single left to right pass where

  - Elements processed in real time

  - Can't store the full data. => use minimal amount of storage while maintaining working "summary"

# Counting distinct elements

32, 12, 14, 32, 7, 12, 32, 7, 32, 12, 4,

$x_1$  $x_2$  $x_3$  $x_4$  . . .

output  m = 5

Applications:
- IP packet streams: How many distinct IP addresses or IP flows (source+destination IP, port, protocol)

    - Anomaly detection, traffic monitoring
- Search: How many distinct search queries on Google on a certain topic yesterday
- Web services: how many distinct users (cookies) searched/browsed a certain term/item

    - Advertising, marketing trends, etc.

# Counting distinct elements

32, 12, 14, 32, 7, 12, 32, 7, 32, 12, 4,

- Want to compute number of **distinct** keys in the stream.

- *How to do this without storing all the elements?*

- *Yet another super cool application of probability (and hashing)*

# A naive solution, counting!

Store the **m** **distinct** user IDs
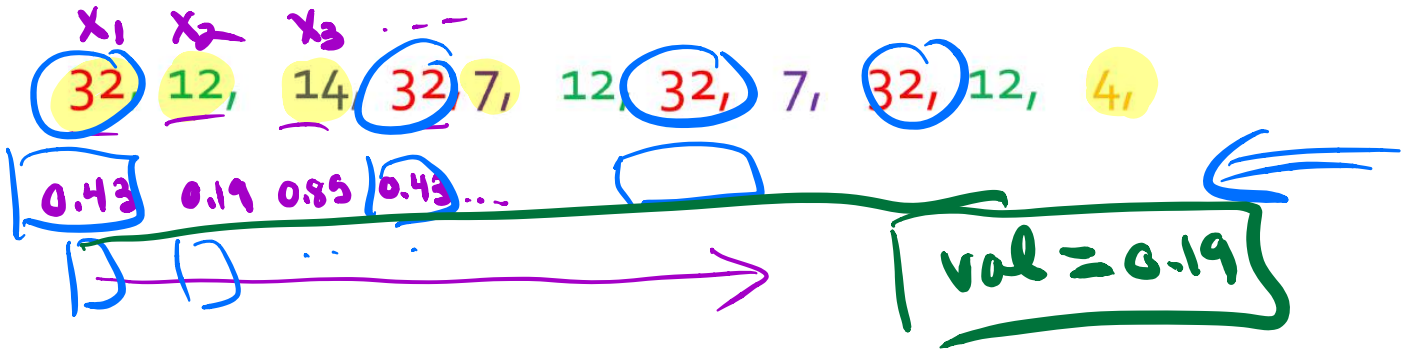in a hash table.

Space requirement: O(**m**)

# CONSIDERING THE NUMBER OF USERS OF YOUTUBE, AND THE NUMBER OF VIDEOS ON YOUTUBE, THIS IS NOT FEASIBLE.

Consider a hash function $h : \mathcal{U} \to [0, 1]$
For distinct values in $\mathcal{U}$ , the
function maps to iid (independent and
identically distributed) Unif(0,1)
random numbers.

Note that, if you were to feed in two
equivalent elements, the function
returns the **same** number.

$h(32) = 0.43$
$h(12) = 0.19$
$h(14) = 0.85$
$h(7) = 0.61$

$X_1 \quad X_2 \quad X_3 \quad \cdots$

32, 12, 14, 32, 7, 12, 32, 7, 32, 12, 4,

0.43  0.19  0.85  0.43 ...

val = 0.19

Track minimum hash value seen so far

$$\min\left(h(x_1), \ldots, h(x_N)\right)$$

In this ex.
0.19

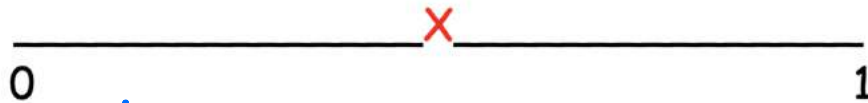$$= \min\left(U_1, U_2, \ldots, U_m\right)$$

$U_i \sim U(0,1)$

$m$ # distinct elts

# MIN OF IID UNIFORMS

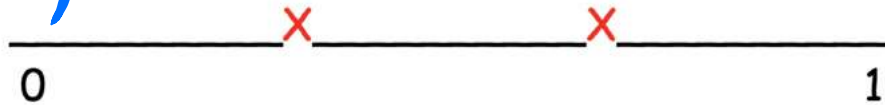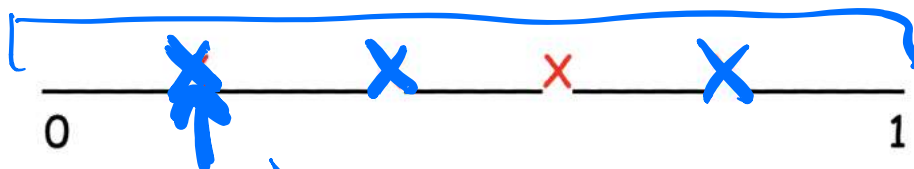If $Y_1, \ldots, Y_m$ are iid $Unif(0,1)$, where do we "expect" the points to end up?

$E(Y_1)$

$m = 1$

$\qquad \qquad \qquad \times$

0 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ 1

$m = 2$

$E(\min(Y_1, Y_2)) = \frac{1}{3}$

$\qquad\qquad \times \qquad\qquad\qquad \times$

0 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ 1

$m = 4$

$\qquad \times \quad \times \quad \times \quad \times$

0 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ 1

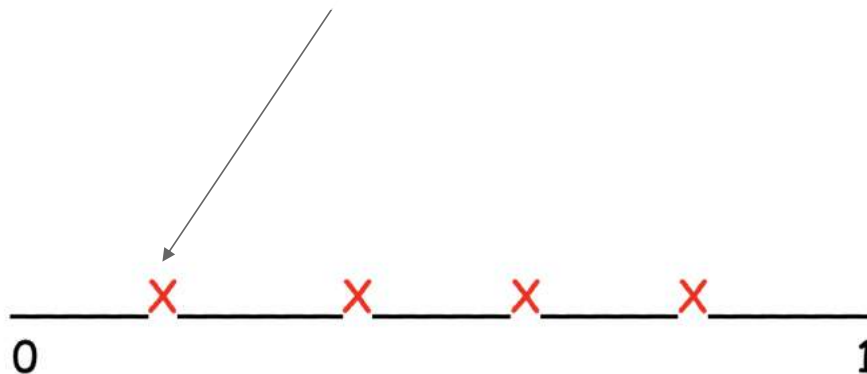$E(\min(Y_1, Y_2, Y_3, Y_4)) = \frac{1}{5}$

# Min of IID Uniforms

If $Y_1, \ldots, Y_m$ are iid $Unif(0,1)$, where do we "expect" the points to end up?

$$E[\min\{Y_1, \ldots, Y_4\}] = \frac{1}{4+1} = \frac{1}{5}$$

$m = 4$

$$\underline{\hspace{1.2cm}}\overset{\textcolor{red}{\times}}{\hspace{1.4cm}}\underline{\hspace{1.4cm}}\overset{\textcolor{red}{\times}}{\hspace{0.8cm}}\underline{\hspace{1.4cm}}\overset{\textcolor{red}{\times}}{\hspace{0.8cm}}\underline{\hspace{1.4cm}}\overset{\textcolor{red}{\times}}{\hspace{0.8cm}}\underline{\hspace{2cm}}$$

0                                               1

# MIN OF IID UNIFORMS

If $Y_1, \ldots, Y_m$ are iid $Unif(0,1)$, where do we "expect" the points to end up?

$$E\left[\min\left(Y_1, \ldots, Y_m\right)\right] = \frac{1}{m+1}$$

$$Y: iid \; U(0,1)$$

- proved in section mbook
- section worksheets
- CC today

# A super duper clever idea



$$E\left[\min\left(h(x_1), \ldots, h(x_N)\right)\right] = \frac{1}{m+1}$$

$m$   # distinct elts.

val

$$E(val) = \frac{1}{m+1}$$

$$= \quad m+1 = \frac{1}{E(val)} \quad \Rightarrow \quad m = \frac{1}{E(val)} - 1$$

What if   $val \approx E(val)$

estimate $m$ to be $\frac{1}{val} - 1$

rounded to nearest int (floor)

$$h(32) = 0.43$$
$$h(5) = 0.19$$
$$h(17) = 0.85$$
$$h(4) = 0.61$$

$x_1 \quad x_2 \quad x_3 \qquad\qquad\qquad\qquad x_N$

| 32 | 5 | 17 | 32 | **14** | 5 | 32 | 32 | 17 |

$$\min_{1 \le i \le N} h(x_i) = 0.19 \qquad (\text{val})$$

$$\hat{m} = \text{round}\left(\frac{1}{0.19} - 1\right)$$

estimate
of # distinct elts.

4.26 ⟸

# The Distinct Elements Algorithm

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
  val ← ∞
**function** UPDATE(x)
  val ← min {val, hash(x)}
**function** ESTIMATE()
  **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**    ▷ Loop through all stream elements
  update($x_i$)    ▷ Update our single float variable
**return** estimate()    ▷ An estimate for $n$, the number of distinct elements.

val tracks
$$\min \left( h(x_1), \ldots, h(x_N) \right)$$

# Distinct Elements Example

Stream:   13,    25,    19,    25,    19,    19

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    val $\leftarrow \infty$
**function** UPDATE(x)
    val $\leftarrow$ min {val, hash(x)}
**function** ESTIMATE()
    **return** round $\left(\frac{1}{\text{val}} - 1\right)$

**for** $i = 1, \ldots, N$: **do**        ▷ Loop through all stream elements
    update($x_i$)        ▷ Update our single float variable
**return** estimate()        ▷ An estimate for $n$, the number of distinct elements.

`val = infty`

# Distinct Elements Example

Stream:  13,    25,    19,    25,    19,    19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    val ← ∞
**function** UPDATE(x)
    val ← min {val, hash(x)}
**function** ESTIMATE()
    **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**         ▷ Loop through all stream elements
    update($x_i$)            ▷ Update our single float variable
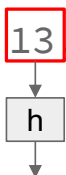**return** estimate()     ▷ An estimate for $n$, the number of distinct elements.

`val = infty`

# Distinct Elements Example

Stream: 13, 25, 19, 25, 19, 19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    val ← ∞
**function** UPDATE(x)
    val ← min {val, hash(x)}
**function** ESTIMATE()
    **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**          ▷ Loop through all stream elements
    update($x_i$)               ▷ Update our single float variable
**return** estimate()          ▷ An estimate for $n$, the number of distinct elements.

val = 0.51

# Distinct Elements Example

Stream:  13,  25,  19,  25,  19,  19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

---

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
   val ← ∞
**function** UPDATE(x)
   val ← min {val, hash(x)}
**function** ESTIMATE()
   **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**                                      ▷ Loop through all stream elements
   update($x_i$)                                                    ▷ Update our single float variable
**return** estimate()                        ▷ An estimate for $n$, the number of distinct elements.

val = 0.26

# Distinct Elements Example

Stream:  13,   25,   19,   25,   19,   19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    val ← ∞
**function** UPDATE(x)
    val ← min {val, hash(x)}
**function** ESTIMATE()
    **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**          ▷ Loop through all stream elements
    update($x_i$)                  ▷ Update our single float variable
**return** estimate()          ▷ An estimate for $n$, the number of distinct elements.

val = 0.26

# Distinct Elements Example

Stream:  13,   25,   19,   25,   19,   19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    val $\leftarrow \infty$
**function** UPDATE(x)
    val $\leftarrow$ min {val, hash(x)}
**function** ESTIMATE()
    **return** round $\left( \frac{1}{val} - 1 \right)$

**for** $i = 1, \ldots, N$: **do**        ▷ Loop through all stream elements
    update($x_i$)        ▷ Update our single float variable
**return** estimate()    ▷ An estimate for $n$, the number of distinct elements.

**val = 0.26**

# Distinct Elements Example

Stream:  13,   25,   19,   25,   19,   19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    val ← ∞
**function** UPDATE(x)
    val ← min {val, hash(x)}
**function** ESTIMATE()
    **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**     ▷ Loop through all stream elements
    update($x_i$)     ▷ Update our single float variable
**return** estimate()     ▷ An estimate for $n$, the number of distinct elements.

val = 0.26

# Distinct Elements Example

Stream:  13,   25,   19,   25,   19,   19

$$\boxed{h}$$

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
   val ← ∞
**function** UPDATE(x)
   val ← min {val, hash(x)}
**function** ESTIMATE()
   **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**          ▷ Loop through all stream elements
   update($x_i$)          ▷ Update our single float variable
**return** estimate()          ▷ An estimate for $n$, the number of distinct elements.

*(handwritten annotations)*

$val_1, \ldots, val_i$

$val_i := \min(val_i, h_i(x))$

using i2

$\hat{val} = \frac{1}{k} \sum_{i=1}^{k} val_i$

val = 0.26

# DISTINCT ELEMENTS EXAMPLE

h: maps $V$ to **randomly** in $(0,1)$

Stream:  13,   25,   19,   25,   19,   19

h

Hashes: 0.51, 0.26, 0.79, 0.26, 0.79, 0.79

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
  val ← ∞
**function** UPDATE(x)
  val ← min {val, hash(x)}
**function** ESTIMATE()
  **return** round $\left(\frac{1}{val} - 1\right)$
**for** $i = 1, \ldots, N$: **do**          ▷ Loop through all stream elements
  update($x_i$)                      ▷ Update our single float variable
**return** estimate()          ▷ An estimate for $n$, the number of distinct elements.

**val = 0.26**

Return
round(1/0.26 − 1) =
round(2.846) =
3

# Diy: Distinct Elements Example II

Stream: 11, 34, 89, 11, 89, 23

Hashes: 0.5, 0.21, 0.94, 0.5, 0.94, 0.1

**Algorithm 2** Distinct Elements Operations

**function** INITIALIZE()
    $val \leftarrow \infty$
**function** UPDATE(x)
    $val \leftarrow \min\{val, hash(x)\}$
**function** ESTIMATE()
    **return** round $\left(\frac{1}{val} - 1\right)$

**for** $i = 1, \ldots, N$: **do**          ▷ Loop through all stream elements
    update($x_i$)          ▷ Update our single float variable
**return** estimate()      ▷ An estimate for $n$, the number of distinct elements.

val = 0.1

Return= 9

# SUMMARY SO FAR

$$h: \underset{\text{universe}}{U} \longrightarrow (0,1)$$

track $\quad val := \min_{1 \leq i \leq N} h(x_i)$

output $\quad \text{round}\left(\frac{1}{val} - 1\right)$

Assume $h$ is
really great
in that $\quad \forall x \leftarrow U$
$h(x) \sim U(0,1)$
indep.

$$h': U \to \{0, 1, \ldots, H-1\}$$
$$h(x) = \frac{h'(x)}{H}$$

$m$ distinct elts $x_1, \ldots, x_N$

PROBLEM   Is val likely to close to $E(val)$?

$m-\sigma$   $m+\sigma$   $0$   $1$

$$E(val) = \frac{1}{m+1}$$

$$Var(val) \approx \frac{1}{(m+1)^2}$$

$$\sigma(val) \approx \frac{1}{m+1} \Longleftarrow$$

$$\hat{m} = round\left(\left(\frac{1}{val}\right)^{-1}\right)$$

for many random vars including val, likely to be $> 0.999 \sigma$ away from mean

# HOW CAN WE REDUCE THE VARIANCE?

## Repetitions

Let $h_1, \ldots, h_k$ be indep hash fns.

$$h_i(x_j) \sim U(0,1) \qquad [\text{assume all of these indep}]$$

track

$$val_1 = \min\left( h_1(x_1) \ldots , h_1(x_N) \right) = \min\left( U_1^1, \ldots, U_m^1 \right)$$

$$val_2 = \min\left( h_2(x_1) \ldots h_2(x_N) \right) = \min\left( U_1^2, \ldots, U_m^2 \right)$$

$$\vdots$$

$$val_k = \min\left( h_k(x_1) \ldots , h_k(x_N) \right) = \min\left( U_1^k \ldots U_m^k \right)$$

$$\overline{val} = \frac{1}{K}\left(val_1 + val_2 + \ldots + val_K\right)$$

$$E(val_i) = \frac{1}{m+1}$$

$$E(\widehat{val}) = E\left(\frac{1}{K}\left(val_1 + val_2 + \ldots + val_K\right)\right)$$

$$= \frac{1}{K}\left[E(val_1) + E(val_2) + \ldots E(val_K)\right]$$
$$\underbrace{\phantom{E(val_1)}}_{\frac{1}{m+1}} \quad \underbrace{\phantom{E(val_2)}}_{\frac{1}{m+1}}$$

$$= \frac{1}{K} \quad \frac{K}{m+1} = \frac{1}{m+1}$$

$$Var(\widehat{val}) = Var\left(\boxed{\frac{1}{K}}\left(val_1 + \ldots + val_K\right)\right)$$

$Var(aX)$
$= a^2 Var(X)$

$$= \frac{1}{K^2} Var\left(val_1 + \ldots + val_K\right)$$

$$\overset{Indep}{=} \frac{1}{K^2}\left[Var(val_1) + \ldots + Var(val_K)\right]$$
$$\underbrace{\phantom{Var(val_1)}}_{=\frac{1}{(m+1)^2}}$$

**$E(\widehat{val})$**

a) $\dfrac{1}{m+1}$

b) $\dfrac{1}{K(m+1)}$

c) $\dfrac{K}{(m+1)}$

d) I don't know

**$Var(X)$**

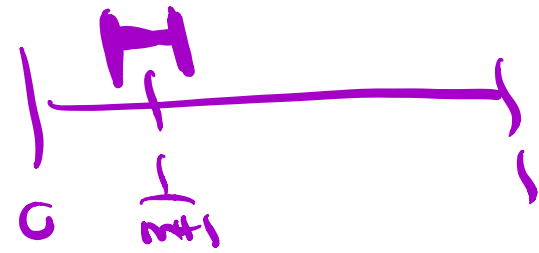a) $\approx \dfrac{1}{(m+1)^2}$

b) $\approx \dfrac{1}{K^2(m+1)^2}$

c) $\approx \dfrac{1}{K(m+1)^2}$

d) I don't know?

$$= \frac{1}{k^2} \quad k \cdot \frac{1}{(m+1)^2} = \frac{1}{k(m+1)^2}$$

$$Var(\widehat{val_i}) \simeq \frac{1}{(m+1)^2}$$

$$\sigma\left(\widehat{val}\right) = \frac{1}{\sqrt{k}(m+1)}$$



$0 \qquad \frac{1}{m+1}$

$$\widehat{val_i} = \min\left(h(x_1), \dots, h(x_N)\right)$$

$x_i = x_j$

$h(x_i) = h(x_j)$

$$= \min\left(U_1, \dots, U_m\right)$$

$x_1 \dots, \quad x_N$

$h(x_{i_1}) \qquad h(x_{i_m})$

m distinct elts

$x_{i_1}, \dots \ x_{i_m}$

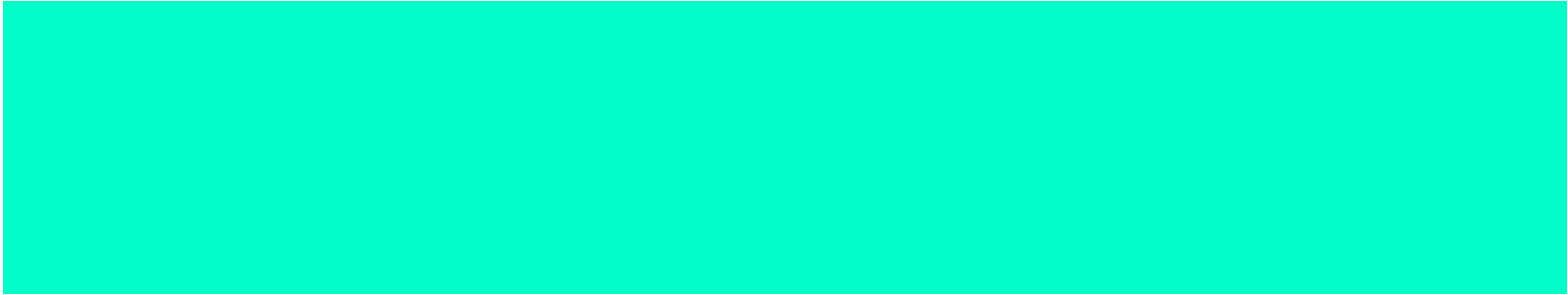$$E(val) = E\left[\min\left(U_1, \dots, U_m\right)\right] = \frac{1}{m+1}$$

# Coding on pset 6

You will use a hash function $h : \mathcal{U} \to [0, 1]$
For distinct values in $\mathcal{U}$, the function
maps to iid (independent and identically
distributed) Unif(0,1) random numbers.

Note that, if you were to feed in two
equivalent elements, the function returns
the **same** number.

We will implement the hash function for
you! Just know that you can consider it an
iid uniform continuous random variables
for each of the values being hashed.

# To do better...

1. we will keep track of K DistElts classes each with its own independent hash function
2. take the mean of our K mins to get a better estimate of the min
3. and then apply the same trick as earlier to give an estimate for the number of distinct elements based on this min that we saw.

# Joint Distributions

Anna Karlin
Most slides by Alex Tsun

# 5.1 Joint Discrete Distributions

# Agenda

- Motivation
- Cartesian Products of Sets
- Joint PMFs and Expectation
- Marginal PMFs

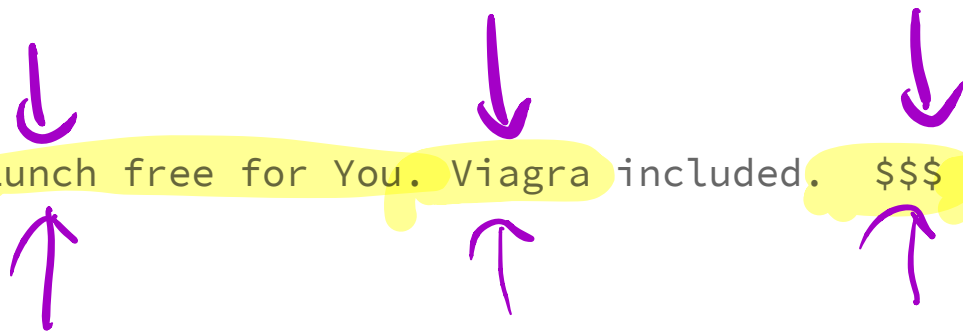# Naive Bayes Classifier - What we Calculate

$$\mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"}) = \frac{\mathbb{P}(\text{"You buy Viagra!"} \mid \text{spam})\,\mathbb{P}(\text{spam})}{\mathbb{P}(\text{"You buy Viagra!"})}$$

$$= \frac{\mathbb{P}(\{\text{"you"},\text{"buy"},\text{"viagra"}\} \mid \text{spam})\,\mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"},\text{"buy"},\text{"viagra"}\} \mid \text{spam})\,\mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"},\text{"buy"},\text{"viagra"}\} \mid \text{ham})\,\mathbb{P}(\text{ham})}$$

[LTP]

# Naive Bayes Classifier – The naive part

$$\mathbb{P}(\{\text{``you''}, \text{``buy''}, \text{``viagra''}\} \mid \text{spam})$$
$$\approx \mathbb{P}(\text{``you''} \mid \text{spam})\mathbb{P}(\text{``buy''} \mid \text{spam})\mathbb{P}(\text{``viagra''} \mid \text{spam})$$

# Why is this Naive?

"!!!Lunch free for You. Viagra included.  $$$ !!!!!"

# Ubiquitous in ML

$$Pr\left(x_1=35, BP=95, \ldots, disease=No\right)$$

- Given "labeled data" $x_1$ $x_2$ $x_3$     $x_5$

| Temp. | BP. | Sore Throat | ... | Colour | diseaseX |
|-------|-----|-------------|-----|--------|----------|
| 35 | 95 | Y | ... | Pale | No |
| 22 | 110 | N | ... | Clear | Yes |
| : | : | | | : | : |
| 10 | 87 | N | ... | Pale | No |

- Learn CLASSIFIER, that can predict label of *NEW* instance

**Learner**

| Temp | BP | Sore-Throat | ... | Color | diseaseX |
|------|-----|-------------|-----|-------|----------|
| 32 | 90 | N | ... | Pale | ? |

**Classifier**

| diseaseX |
|----------|
| **No** |

39