

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, \dots, 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,

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 n **distinct** user IDs
in a hash table.

Space requirement: $O(n)$



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

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

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

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

MIN OF IID UNIFORMS

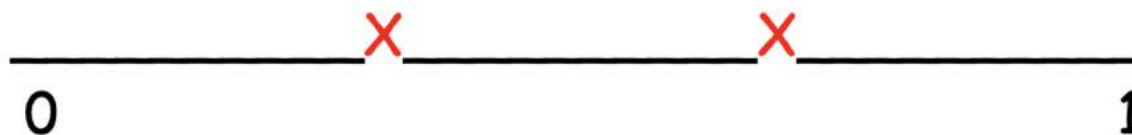


If Y_1, \dots, Y_m are iid $Unif(0,1)$, where do we "expect" the points to end up?

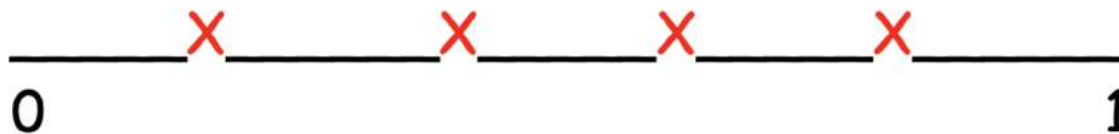
$m = 1$



$m = 2$



$m = 4$



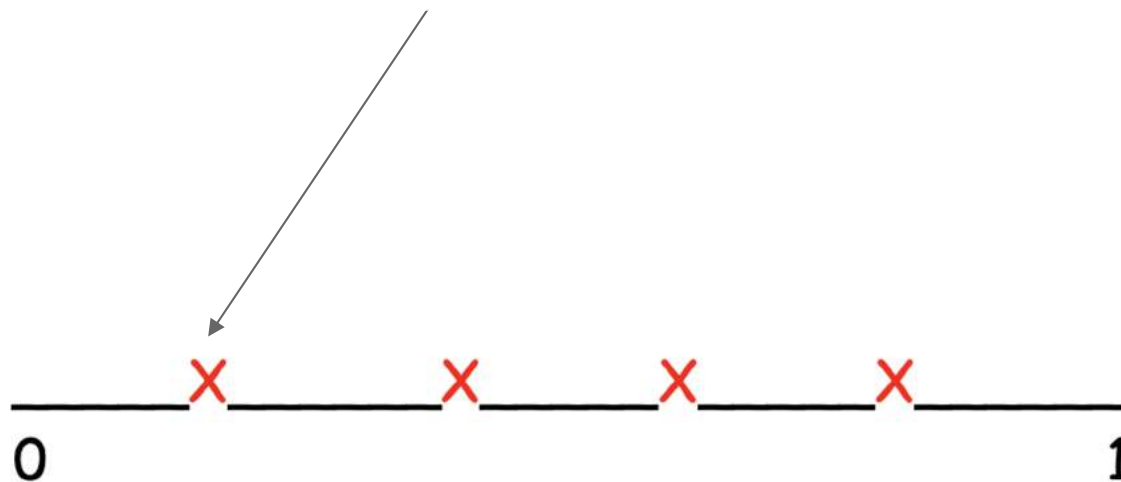
MIN OF IID UNIFORMS



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

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

$m = 4$



MIN OF IID UNIFORMS



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

A SUPER DUPER CLEVER IDEA



32 5 17 32 14 5 32 32 17



THE DISTINCT ELEMENTS ALGORITHM

Algorithm 2 Distinct Elements Operations

function INITIALIZE()

$\text{val} \leftarrow \infty$

function UPDATE(x)

$\text{val} \leftarrow \min \{ \text{val}, \text{hash}(x) \}$

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

$\text{update}(x_i)$

return $\text{estimate}()$

▷ Loop through all stream elements

▷ Update our single float variable

▷ An estimate for n , the number of distinct elements.

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 ∞

function UPDATE(x)

val \leftarrow min {val, hash(x)}

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = infity

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 ∞

function UPDATE(x)

val \leftarrow min {val, hash(x)}

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

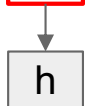
▸ An estimate for n , the number of distinct elements.

val = infty

DISTINCT ELEMENTS EXAMPLE

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

13



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

0.51

Algorithm 2 Distinct Elements Operations

function INITIALIZE()

val $\leftarrow \infty$

function UPDATE(x)

val $\leftarrow \min \{ \text{val}, \text{hash}(x) \}$

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = 0.51

DISTINCT ELEMENTS EXAMPLE

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

25

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 $\leftarrow \min \{ \text{val}, \text{hash}(x) \}$

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = 0.26

DISTINCT ELEMENTS EXAMPLE

Stream: 13, 25, 19, 25, 19, 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 ∞

function UPDATE(x)

val \leftarrow min {val, hash(x)}

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = 0.26

DISTINCT ELEMENTS EXAMPLE

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

25

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 $\leftarrow \min \{ \text{val}, \text{hash}(x) \}$

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = 0.26

DISTINCT ELEMENTS EXAMPLE

Stream: 13, 25, 19, 25, 19, 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 ∞

function UPDATE(x)

val \leftarrow min {val, hash(x)}

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ 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 $\leftarrow \min \{ \text{val}, \text{hash}(x) \}$

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

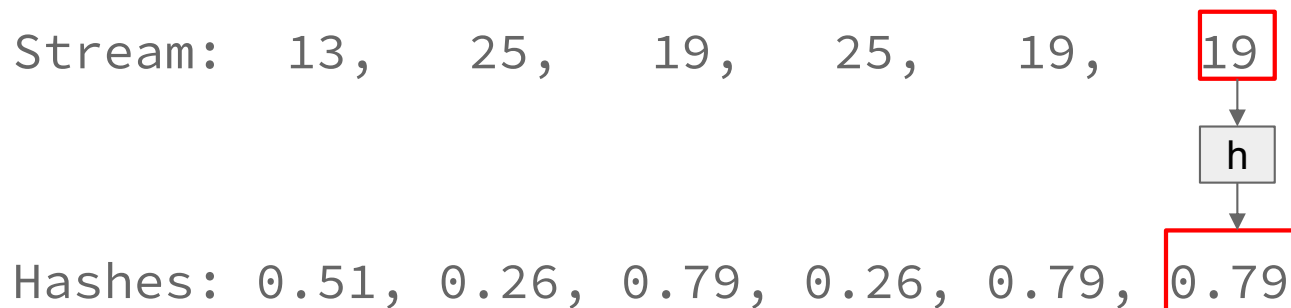
▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = 0.26

DISTINCT ELEMENTS EXAMPLE



Algorithm 2 Distinct Elements Operations

function INITIALIZE()

val $\leftarrow \infty$

function UPDATE(x)

val $\leftarrow \min \{ \text{val}, \text{hash}(x) \}$

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ 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 ∞

function UPDATE(x)

val \leftarrow min {val, hash(x)}

function ESTIMATE()

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

for $i = 1, \dots, N$: **do**

update(x_i)

return estimate()

▸ Loop through all stream elements

▸ Update our single float variable

▸ An estimate for n , the number of distinct elements.

val = 0.1

Return= 9

SUMMARY SO FAR

PROBLEM

HOW CAN WE REDUCE THE VARIANCE?



CODING ON PSET 6

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

$$\begin{aligned} \mathbb{P}(\text{spam} \mid \text{"You buy Viagra!"}) &= \frac{\mathbb{P}(\text{"You buy Viagra!"} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\text{"You buy Viagra!"})} \\ &= \frac{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam})}{\mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{spam}) \mathbb{P}(\text{spam}) + \mathbb{P}(\{\text{"you"}, \text{"buy"}, \text{"viagra"}\} \mid \text{ham}) \mathbb{P}(\text{ham})} \quad [\text{LTP}] \end{aligned}$$

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. \$\$\$!!!!!”

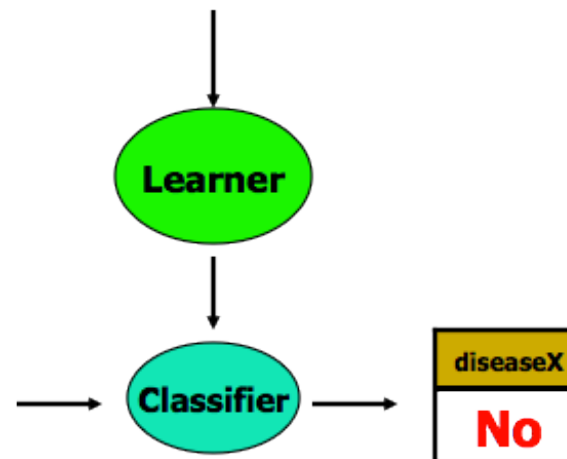
UBIQUITOUS IN ML

- Given “labeled data”

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

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



AGENDA

- CARTESIAN PRODUCTS OF SETS
- JOINT PMFS AND EXPECTATION
- MARGINAL PMFS

CARTESIAN PRODUCT OF SETS

Cartesian Product: Let A, B be sets. The Cartesian product of A and B is denoted

$$A \times B = \{(a, b): a \in A, b \in B\}$$

A small example:

$$\{1, 2, 3\} \times \{4, 5\} = \{(1, 4), (1, 5), (2, 4), (2, 5), (3, 4), (3, 5)\}$$

Another example: The xy-plane (2D space) is denoted

$$\mathbb{R}^2 = \mathbb{R} \times \mathbb{R} = \{(x, y): x \in \mathbb{R}, y \in \mathbb{R}\}$$

If A, B are finite sets, then $|A \times B| = |A| \cdot |B|$ by the product rule of counting.

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value
of the red die. Specify

$$\Omega_X = \{1,2,3,4\}$$

$$\Omega_Y = \{1,2,3,4\}$$

$$\Omega_{X,Y} = \Omega_X \times \Omega_Y$$

Specify the joint PMF $p_{X,Y}(x,y) = P(X = x, Y = y)$
for $x, y \in \Omega_{X,Y}$.

$X \setminus Y$	1	2	3	4
1				
2				
3				
4				

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value
of the red die. Specify

$$\Omega_X = \{1,2,3,4\}$$

$$\Omega_Y = \{1,2,3,4\}$$

$$\Omega_{X,Y} = \Omega_X \times \Omega_Y$$

Specify the joint PMF $p_{X,Y}(x,y) = P(X = x, Y = y)$
for $x, y \in \Omega_{X,Y}$.

$$p_{X,Y}(x,y) = \begin{cases} 1/16, & x, y \in \Omega_{X,Y} \\ 0, & \text{otherwise} \end{cases}$$

$X \setminus Y$	1	2	3	4
1	1/16	1/16	1/16	1/16
2	1/16	1/16	1/16	1/16
3	1/16	1/16	1/16	1/16
4	1/16	1/16	1/16	1/16

JOINT PMFS AND EXPECTATION

Joint PMFs: Let X, Y be discrete random variables. The joint PMF of X and Y is

$$p_{X,Y}(a, b) = P(X = a, Y = b)$$

The joint range is

$$\Omega_{X,Y} = \{(c, d) : p_{X,Y}(c, d) > 0\} \subseteq \Omega_X \times \Omega_Y$$

Note that

$$\sum_{(s,t) \in \Omega_{X,Y}} p_{X,Y}(s, t) = 1$$

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value of the red die. Let $U = \min \{X, Y\}$ and $V = \max \{X, Y\}$.

$$\Omega_U = \{1, 2, 3, 4\}$$

$$\Omega_V = \{1, 2, 3, 4\}$$

$$\Omega_{U,V} = \{(u, v) \in \Omega_U \times \Omega_V : u \leq v\} \neq \Omega_U \times \Omega_V$$

$U \setminus V$	1	2	3	4
1				
2				
3				
4				

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value of the red die. Let $U = \min \{X, Y\}$ and $V = \max \{X, Y\}$.

$$\Omega_U = \{1, 2, 3, 4\}$$

$$\Omega_V = \{1, 2, 3, 4\}$$

$$\Omega_{U,V} = \{(u, v) \in \Omega_U \times \Omega_V: u \leq v\} \neq \Omega_U \times \Omega_V$$

Specify the joint PMF $p_{U,V}(u, v) = P(U = u, V = v)$
for $u, v \in \Omega_{U,V}$.

UV	1	2	3	4
1				
2				
3				
4				

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value of the red die. Let $U = \min \{X, Y\}$ and $V = \max \{X, Y\}$.

$$\Omega_U = \{1,2,3,4\}$$

$$\Omega_V = \{1,2,3,4\}$$

$$\Omega_{U,V} = \{(u,v) \in \Omega_U \times \Omega_V: u \leq v\} \neq \Omega_U \times \Omega_V$$

Specify the joint PMF $p_{U,V}(u,v) = P(U = u, V = v)$
for $u, v \in \Omega_{U,V}$.

$$p_{U,V}(u,v) = \begin{cases} 2/16, & u, v \in \Omega_U \times \Omega_V, & v > u \\ 1/16, & u, v \in \Omega_U \times \Omega_V, & v = u \\ 0, & \text{otherwise} \end{cases}$$

$U \setminus V$	1	2	3	4
1	1/16	2/16	2/16	2/16
2	0	1/16	2/16	2/16
3	0	0	1/16	2/16
4	0	0	0	1/16

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value of the red die. Let $U = \min \{X, Y\}$ and $V = \max \{X, Y\}$.

What is $p_U(u)$ for $u \in \Omega_U$?

$$p_U(u) = \begin{cases} u = 1 \\ u = 2 \\ u = 3 \\ u = 4 \end{cases}$$

$U \setminus V$	1	2	3	4
1	1/16	2/16	2/16	2/16
2	0	1/16	2/16	2/16
3	0	0	1/16	2/16
4	0	0	0	1/16

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value
of the red die. Let $U = \min \{X, Y\}$ and $V = \max \{X, Y\}$.

What is $p_U(u)$ for $u \in \Omega_U$?

$$p_U(u) = \begin{cases} 7/16, & u = 1 \\ 5/16, & u = 2 \\ 3/16, & u = 3 \\ 1/16, & u = 4 \end{cases}$$

$U \setminus V$	1	2	3	4
1	1/16	2/16	2/16	2/16
2	0	1/16	2/16	2/16
3	0	0	1/16	2/16
4	0	0	0	1/16

EXAMPLE: WEIRD DICE AGAIN



Suppose I roll two fair 4-sided die independently.
Let X be the value of the blue die, and Y the value of the red die. Let $U = \min \{X, Y\}$ and $V = \max \{X, Y\}$.

UV	1	2	3	4
1	1/16	2/16	2/16	2/16
2	0	1/16	2/16	2/16
3	0	0	1/16	2/16
4	0	0	0	1/16

JOINT PMFS AND EXPECTATION

Joint PMFs: Let X, Y be discrete random variables. The joint PMF of X and Y is

$$p_{X,Y}(a, b) = P(X = a, Y = b)$$

The joint range is

$$\Omega_{X,Y} = \{(c, d): p_{X,Y}(c, d) > 0\} \subseteq \Omega_X \times \Omega_Y$$

Note that

$$\sum_{(s,t) \in \Omega_{X,Y}} p_{X,Y}(s, t) = 1$$

If $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ is a function, then

$$E[g(X, Y)] = \sum_x \sum_y g(x, y) p_{X,Y}(x, y)$$

MARGINAL PMFS

Marginal PMFs: Let X, Y be discrete random variables. The marginal PMF of X is

$$p_X(a) = \sum_{b \in \Omega_Y} p_{X,Y}(a, b)$$

MARGINAL PMFS

Marginal PMFs: Let X, Y be discrete random variables. The marginal PMF of X is

$$p_X(a) = \sum_{b \in \Omega_Y} p_{X,Y}(a, b)$$

Similarly, the marginal PMF of Y is

$$p_Y(d) = \sum_{c \in \Omega_X} p_{X,Y}(c, d)$$

MARGINAL PMFS

Marginal PMFs: Let X, Y be discrete random variables. The marginal PMF of X is

$$p_X(a) = \sum_{b \in \Omega_Y} p_{X,Y}(a, b)$$

Similarly, the marginal PMF of Y is

$$p_Y(d) = \sum_{c \in \Omega_X} p_{X,Y}(c, d)$$

(Extension) If Z is also a discrete random variable, then the marginal PMF of Z is

$$p_Z(z) = \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} p_{X,Y,Z}(x, y, z)$$

INDEPENDENCE

Independence (DRVs): Discrete random variables X, Y are independent, written $X \perp Y$, if for all $x \in \Omega_X$ and $y \in \Omega_Y$,

$$p_{X,Y}(x, y) = p_X(x)p_Y(y)$$

Recall $\Omega_{X,Y} = \{(x, y): p_{X,Y}(x, y) > 0\} \subseteq \Omega_X \times \Omega_Y$. A necessary but not sufficient condition for independence is that $\Omega_{X,Y} = \Omega_X \times \Omega_Y$. That is, if $\Omega_{X,Y} \neq \Omega_X \times \Omega_Y$, then X and Y cannot be independent, but if $\Omega_{X,Y} = \Omega_X \times \Omega_Y$, then we have to check the condition.

This is because if there is some $(a, b) \in \Omega_X \times \Omega_Y$ but not in $\Omega_{X,Y}$, then $p_{X,Y}(a, b) = 0$ but $p_X(a) > 0$ and $p_Y(b) > 0$, violating independence.

VARIANCE ADDS FOR INDEPENDENT RVS

If X, Y are independent random variables $X \perp Y$, then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

This property relies on the fact that they are independent, whereas linearity of expectation **always** holds, regardless. If $a, b, c \in \mathbb{R}$ are scalars, then

$$\text{Var}(aX + bY + c) = a^2\text{Var}(X) + b^2\text{Var}(Y)$$

If $X \perp Y$, then $E[XY] = E[X]E[Y]$.

