# CS294-248 Special Topics in Database Theory
## Unit 4: AGM Bound, WCOJ

Dan Suciu

University of Washington

## Outline

- Today: the AGM bound. This is a mathematical formula that gives us $AGM(Q, \boldsymbol{D}) \overset{\text{def}}{=} \max_{\boldsymbol{D} \models \text{statistics}} |Q(\boldsymbol{D})|$.

- Thursday: Worst Case Optimal Join. This is an algorithm that computes $Q(\boldsymbol{D})$ in time $\tilde{O}(AGM(Q, \boldsymbol{D}))$.

# Background on Cardinality Estimation

# Cardinality Estimation 101 (1/3)

Given:

- Statistics on the input relations $R_1, R_2, \ldots$
- A full conjunctive query $Q$

"Estimate":

- The size $|Q(\boldsymbol{D})|$.

Numerous applications: query optimization, memory provisioning, data partitioning.

# Cardinality Estimation 101 (2/3)

Bottom-up on the query plan:

- Selection $\sigma_p(R)$: assume independence:

$$\boxed{|\sigma_p(R)| \approx \theta_p \cdot |R|}$$

$$\boxed{\theta_{p_1 \wedge p_2} \approx \theta_{p_1} \cdot \theta_{p_2}}$$

Histograms, multidimensional histograms.

- Join $J(A, B, C) = R(A, B) \wedge S(B, C)$: assume preservation of values
  - $|J| \approx |R| \cdot \operatorname{avg}(\deg_S(C|B)) = \frac{|R| \cdot |S|}{|\operatorname{Dom}(S.B)|}$.
  - $|J| \approx |S| \cdot \operatorname{avg}(\deg_R(A|B)) = \frac{|R| \cdot |S|}{|\operatorname{Dom}(R.B)|}$.
  - Heuristic: take the *minimum*: $|J| \approx \frac{|R| \cdot |S|}{\max(|\operatorname{Dom}(R.B)| \cdot |\operatorname{Dom}(S.B)|)}$

# Cardinality Estimation 101 (2/3)

Bottom-up on the query plan:

- Selection $\sigma_p(R)$: assume independence:

$$\boxed{|\sigma_p(R)| \approx \theta_p \cdot |R|}$$

$$\boxed{\theta_{p_1 \wedge p_2} \approx \theta_{p_1} \cdot \theta_{p_2}}$$

  Histograms, multidimensional histograms.

- Join $J(A, B, C) = R(A, B) \wedge S(B, C)$: assume preservation of values
  - $|J| \approx |R| \cdot \text{avg}(\deg_S(C|B)) = \frac{|R| \cdot |S|}{|\text{Dom}(S.B)|}$.
  - $|J| \approx |S| \cdot \text{avg}(\deg_R(A|B)) = \frac{|R| \cdot |S|}{|\text{Dom}(R.B)|}$.
  - Heuristic: take the *minimum*: $|J| \approx \frac{|R| \cdot |S|}{\max(|\text{Dom}(R.B)| \cdot |\text{Dom}(S.B)|)}$

# Cardinality Estimation 101 (2/3)

Bottom-up on the query plan:

- Selection $\sigma_p(R)$: assume independence:

$$\boxed{|\sigma_p(R)| \approx \theta_p \cdot |R|} \qquad\qquad \boxed{\theta_{p_1 \wedge p_2} \approx \theta_{p_1} \cdot \theta_{p_2}}$$

  Histograms, multidimensional histograms.

- Join $J(A, B, C) = R(A, B) \wedge S(B, C)$: assume preservation of values
  - $|J| \approx |R| \cdot \mathrm{avg}(\deg_S(C|B)) = \frac{|R| \cdot |S|}{|\mathrm{Dom}(S.B)|}$.
  - $|J| \approx |S| \cdot \mathrm{avg}(\deg_R(A|B)) = \frac{|R| \cdot |S|}{|\mathrm{Dom}(R.B)|}$.
  - Heuristic: take the *minimum*: $|J| \approx \frac{|R| \cdot |S|}{\max(|\mathrm{Dom}(R.B)| \cdot |\mathrm{Dom}(S.B)|)}$

# Cardinality Estimation 101 (2/3)

Bottom-up on the query plan:

- Selection $\sigma_p(R)$: assume independence:

$$\boxed{|\sigma_p(R)| \approx \theta_p \cdot |R|} \qquad\qquad \boxed{\theta_{p_1 \wedge p_2} \approx \theta_{p_1} \cdot \theta_{p_2}}$$

  Histograms, multidimensional histograms.

- Join $J(A, B, C) = R(A, B) \wedge S(B, C)$: assume preservation of values
  - $|J| \approx |R| \cdot \mathrm{avg}(\deg_S(C|B)) = \frac{|R| \cdot |S|}{|\mathrm{Dom}(S.B)|}$.
  - $|J| \approx |S| \cdot \mathrm{avg}(\deg_R(A|B)) = \frac{|R| \cdot |S|}{|\mathrm{Dom}(R.B)|}$.
  - Heuristic: take the <u>minimum</u>: $\boxed{|J| \approx \frac{|R| \cdot |S|}{\max(|\mathrm{Dom}(R.B)| \cdot |\mathrm{Dom}(S.B)|)}}$

# Cardinality Estimation 101 (3/3)

- Notoriously hard to estimate cardinality of complex queries.

- No rigorous definition of the _estimate_: there is no probability space.

- How do we combine multiple sources of information?
  - We had two formulas for the join, why choose min?
  - Given $R(A, B, C)$ and histograms on $A, B, C, AB, AC$,
    how do we estimate $|\sigma_{A=2, B=4, C=6}(R)|$?

Background
oooo

Output Bound
●oo

AGM Bound
ooooo

Proof: Upper Bound
ooooooo

Proof: Lower Bound
ooooo

Extensions
ooooooooooo

Upper Bound on the Output of a Query

## The Output Bound Problem

Given statistics on the input $\boldsymbol{D}$, e.g. cardinalities, # distinct values,

Compute an upper bound $B$:          $|Q(\boldsymbol{D})| \leq B$

Challenge: make $B$ tight.

## Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$. $\qquad\qquad$ $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = ?$

# Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$.  $\qquad\qquad$ $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$

## Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$.                    $\max_{D} |Q(D)| = N^2$
  If $S.Y$ is a key:                                        $\max_{D} |Q(D)| = N$

# Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$. $\qquad\qquad$ $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  If $S.Y$ is a key: $\qquad\qquad\qquad\qquad\qquad$ $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N$

- $Q(X, Y, Z, U) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$. $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = ?$

# Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$.                    $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  If $S.Y$ is a key:                                        $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N$

- $Q(X, Y, Z, U) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$.  $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$

## Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$.  $\qquad \max_D |Q(D)| = N^2$
  If $S.Y$ is a key:  $\qquad\qquad\qquad\qquad\qquad \max_D |Q(D)| = N$

- $Q(X, Y, Z, U) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$.  $\max_D |Q(D)| = N^2$
  Notice the role of an edge cover

# Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$.  $\qquad \max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  If $S.Y$ is a key: $\qquad\qquad\qquad\qquad\qquad\quad \max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N$

- $Q(X, Y, Z, U) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$.  $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  Notice the role of an edge cover

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$.  $\qquad \max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = ?$

Background
oooo

**Output Bound**
ooo●

AGM Bound
ooooo

Proof: Upper Bound
ooooooo

Proof: Lower Bound
ooooo

Extensions
ooooooooooo

# Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$. $\qquad\qquad \max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  If $S.Y$ is a key: $\qquad\qquad\qquad\qquad\qquad \max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N$

- $Q(X, Y, Z, U) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$. $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  Notice the role of an edge cover

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$. $\qquad \max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$

## Simple Examples

Assume $|R| \leq N$, $|S| \leq N$, $|T| \leq N$.

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z)$.          $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  If $S.Y$ is a key:                    $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N$

- $Q(X, Y, Z, U) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$.   $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^2$
  Notice the role of an edge cover

- $Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$.     $\max_{\boldsymbol{D}} |Q(\boldsymbol{D})| = N^{\frac{3}{2}}$
  Here we use a fractional edge cover

Background
○○○○

Output Bound
○○○

AGM Bound
●○○○○

Proof: Upper Bound
○○○○○○○

Proof: Lower Bound
○○○○○

Extensions
○○○○○○○○○○

# AGM Bound: The Statement

Background
oooo

Output Bound
ooo

AGM Bound
o●ooo

Proof: Upper Bound
ooooooo

Proof: Lower Bound
ooooo

Extensions
ooooooooooo

## Fractional Edge Covers

Query $Q$ to hypegraph $G = (V, E)$. $\qquad R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$

## Fractional Edge Covers

Query $Q$ to hypegraph $G = (V, E)$.         $R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$

### Definition

A *fractional edge cover* is $\boldsymbol{w} = (w_e)_{e \in E}$, $w_e \geq 0$:
    $\forall x \in V$, $\sum_{e \in E : x \in e} w_e \geq 1$.

Background
0000

Output Bound
000

AGM Bound
0●000

Proof: Upper Bound
0000000

Proof: Lower Bound
00000

Extensions
0000000000

# Fractional Edge Covers

Query $Q$ to hypegraph $G = (V, E)$.  $\quad\quad\quad R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$

### Definition

A *fractional edge cover* is $\boldsymbol{w} = (w_e)_{e \in E}$, $w_e \geq 0$:
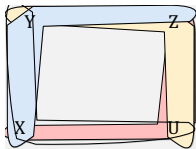$$\forall x \in V, \sum_{e \in E : x \in e} w_e \geq 1.$$
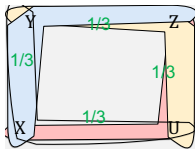
## Examples
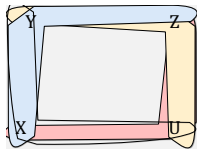
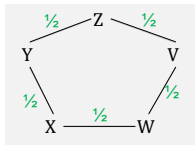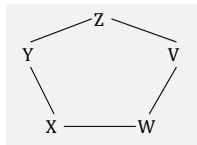What are fractional edge covers?



5-cycle

Loomis-Whitney:

$R(X, Y, Z) \wedge S(Y, Z, U)$
$\wedge T(Z, U, X) \wedge K(U, X, Y)$

## Examples

What are fractional edge covers?





5-cycle

Loomis-Whitney:

$R(X, Y, Z) \wedge S(Y, Z, U)$
$\wedge T(Z, U, X) \wedge K(U, X, Y)$

## Examples

What are fractional edge covers?



5-cycle

Loomis-Whitney:

$R(X, Y, Z) \wedge S(Y, Z, U)$
$\wedge T(Z, U, X) \wedge K(U, X, Y)$

## Examples

What are fractional edge covers?



5-cycle

Loomis-Whitney:

$R(X, Y, Z) \wedge S(Y, Z, U)$
$\wedge T(Z, U, X) \wedge K(U, X, Y)$

## Examples

What are fractional edge covers?



5-cycle

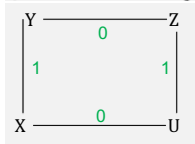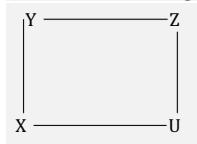Loomis-Whitney:

$R(X, Y, Z) \land S(Y, Z, U)$
$\land T(Z, U, X) \land K(U, X, Y)$

Background
oooo

Output Bound
ooo

AGM Bound
oo●oo

Proof: Upper Bound
ooooooo

Proof: Lower Bound
ooooo

Extensions
oooooooooo

## Examples

What are fractional edge covers?



5-cycle

Loomis-Whitney:

$$R(X, Y, Z) \land S(Y, Z, U)$$
$$\land T(Z, U, X) \land K(U, X, Y)$$

$(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$ is a convex combination of $(1, 0, 1, 0)$ and $(0, 1, 0, 1)$.

## Examples

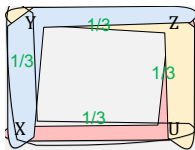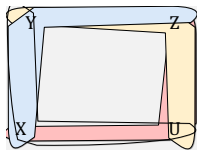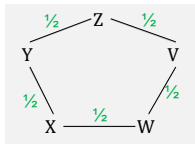What are fractional edge covers?



5-cycle

Loomis-Whitney:

$$R(X, Y, Z) \wedge S(Y, Z, U)$$
$$\wedge T(Z, U, X) \wedge K(U, X, Y)$$

$(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$ is a convex combination of $(1, 0, 1, 0)$ and $(0, 1, 0, 1)$.
Vertex of the edge covering polytope: no convex combination of others.

# The AGM Bound [Atserias et al., 2013]

$Q(\mathbf{X}) = R_1(\mathbf{Y}_1) \wedge \cdots \wedge R_m(\mathbf{Y}_m)$

Theorem (Upper Bound)

*For every fractional edge cover* $\mathbf{w}$: $|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}$

## The AGM Bound [Atserias et al., 2013]

$Q(\mathbf{X}) = R_1(\mathbf{Y}_1) \wedge \cdots \wedge R_m(\mathbf{Y}_m)$

Theorem (Upper Bound)

*For every fractional edge cover $\mathbf{w}$: $|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}$*

Theorem (Lower Bound)

$AGM(Q) \stackrel{def}{=} \min_{\mathbf{w}} |R_1|^{w_1} \cdots |R_m|^{w_m}$ *is "tight".*

## The AGM Bound [Atserias et al., 2013]

$$Q(\boldsymbol{X}) = R_1(\boldsymbol{Y}_1) \wedge \cdots \wedge R_m(\boldsymbol{Y}_m)$$

Theorem (Upper Bound)

*For every fractional edge cover $\boldsymbol{w}$:* $|Q| \le |R_1|^{w_1} \cdots |R_m|^{w_m}$

Theorem (Lower Bound)

$AGM(Q) \stackrel{def}{=} \min_{\boldsymbol{w}} |R_1|^{w_1} \cdots |R_m|^{w_m}$ *is "tight".*

$$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X) \qquad AGM(Q) = \min \begin{pmatrix} (|R| \cdot |S| \cdot |T|)^{1/2} \\ |R| \cdot |S| \\ |R| \cdot |T| \\ |S| \cdot |T| \end{pmatrix}$$

# The AGM Bound [Atserias et al., 2013]

$$Q(\boldsymbol{X}) = R_1(\boldsymbol{Y}_1) \wedge \cdots \wedge R_m(\boldsymbol{Y}_m)$$

### Theorem (Upper Bound)

*For every fractional edge cover $\boldsymbol{w}$: $|Q| \le |R_1|^{w_1} \cdots |R_m|^{w_m}$*

### Theorem (Lower Bound)

$AGM(Q) \overset{def}{=} \min_{\boldsymbol{w}} |R_1|^{w_1} \cdots |R_m|^{w_m}$ *is "tight".*

$$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X) \qquad AGM(Q) = \min \begin{pmatrix} (|R| \cdot |S| \cdot |T|)^{1/2} \\ |R| \cdot |S| \\ |R| \cdot |T| \\ |S| \cdot |T| \end{pmatrix}$$

Minimum over vertices of the edge-covering polytope. WHY?

# Proof Outline

- Proof of the upper bound: information inequalities (a.k.a. entropic inequalities).

- Proof of the lower bound: construct a worst-case database instance by using strong duality of linear optimization.

Background
○○○○

Output Bound
○○○

AGM Bound
○○○○○

Proof: Upper Bound
●○○○○○○

Proof: Lower Bound
○○○○○

Extensions
○○○○○○○○○○○

# Proof of the Upper Bound

Background
0000
Output Bound
000
AGM Bound
00000
Proof: Upper Bound
0●00000
Proof: Lower Bound
00000
Extensions
0000000000

# Entropic Vectors

### Definition

Finite probability space $p : D \to [0, 1]$.　　$X =$ r.v. with outcomes $D$.

The *entropy* of $X$ is:　　　　　$h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$

# Entropic Vectors

## Definition

Finite probability space $p : D \to [0, 1]$.    $X$ = r.v. with outcomes $D$.

The *entropy* of $X$ is:    $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$

$N \stackrel{\text{def}}{=} |D|$:    $0 \leq h(X) \leq \log N$    $h(X) = \log N$ iff $p$ is uniform.

# Entropic Vectors

## Definition

Finite probability space $p : D \to [0, 1]$.     $X =$ r.v. with outcomes $D$.

The *entropy* of $X$ is:     $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$

$N \stackrel{\text{def}}{=} |D|$:     $0 \le h(X) \le \log N$     $h(X) = \log N$ iff $p$ is uniform.

## Definition

R.v. $X_1, \ldots, X_n$. Their *entropic vector* is $\boldsymbol{h} = (h(X_\alpha))_{\alpha \subseteq [n]} \in \mathbb{R}_+^{2^n}$.

# Entropic Vectors

### Definition

Finite probability space $p : D \to [0, 1]$.     $X =$ r.v. with outcomes $D$.

The *entropy* of $X$ is:         $h(X) \overset{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$

$N \overset{\text{def}}{=} |D|$:         $0 \le h(X) \le \log N$         $h(X) = \log N$ iff $p$ is uniform.

### Definition

R.v. $X_1, \ldots, X_n$. Their *entropic vector* is $\boldsymbol{h} = (h(X_\alpha))_{\alpha \subseteq [n]} \in \mathbb{R}_+^{2^n}$.

| X | Y |
|---|---|
| a | p |
| a | q |
| b | q |
| a | m |

# Entropic Vectors

## Definition

Finite probability space $p : D \to [0, 1]$.     $X =$ r.v. with outcomes $D$.

The *entropy* of $X$ is:     $h(X) \stackrel{\text{def}}{=} - \sum_{x \in D} p(x) \log p(x)$

$N \stackrel{\text{def}}{=} |D|$:     $0 \leq h(X) \leq \log N$     $h(X) = \log N$ iff $p$ is uniform.

## Definition

R.v. $X_1, \ldots, X_n$. Their *entropic vector* is $\boldsymbol{h} = (h(X_\alpha))_{\alpha \subseteq [n]} \in \mathbb{R}_+^{2^n}$.

| $X$ | $Y$ | $p$ |
|-----|-----|-----|
| $a$ | $p$ | $1/4$ |
| $a$ | $q$ | $1/4$ |
| $b$ | $q$ | $1/4$ |
| $a$ | $m$ | $1/4$ |

$h(XY) = \log 4$

## Entropic Vectors

### Definition

Finite probability space $p : D \to [0, 1]$.    $X = $ r.v. with outcomes $D$.

The *entropy* of $X$ is:    $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$

$N \stackrel{\text{def}}{=} |D|$:    $0 \leq h(X) \leq \log N$    $h(X) = \log N$ iff $p$ is uniform.

### Definition

R.v. $X_1, \ldots, X_n$. Their *entropic vector* is $\boldsymbol{h} = (h(X_\alpha))_{\alpha \subseteq [n]} \in \mathbb{R}_+^{2^n}$.

| X | Y | p |
|---|---|---|
| a | p | 1/4 |
| a | q | 1/4 |
| b | q | 1/4 |
| a | m | 1/4 |

$h(XY) = \log 4$

| X | p |
|---|---|
| a | 3/4 |
| b | 1/4 |

$h(X) \leq \log 2$

| Y | p |
|---|---|
| p | 1/4 |
| q | 2/4 |
| m | 1/4 |

$h(Y) \leq \log 3$

| $\emptyset$ | p |
|---|---|
|  | 1 |

$h(\emptyset) = 0$

# Shannon Inequalities

### Basic Shannon Inequalities

$$h(\emptyset) = 0$$
$$h(\boldsymbol{U} \cup \boldsymbol{V}) \geq h(\boldsymbol{U}) \qquad\qquad \text{Monotonicity}$$
$$h(\boldsymbol{U}) + h(\boldsymbol{V}) \geq h(\boldsymbol{U} \cup \boldsymbol{V}) + h(\boldsymbol{U} \cap \boldsymbol{V}) \qquad \text{Submodularity}$$

A Shannon inequality is a consequence of these inequalities.

# A Shannon Inequality

### Example

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

## A Shannon Inequality

### Example

$$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$$

$$h(XY) + h(YZ) + h(XZ)$$

## A Shannon Inequality

Example

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

$\underline{h(XY) + h(YZ)} + h(XZ)$

## A Shannon Inequality

### Example

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

$$\underline{h(XY) + h(YZ)} + h(XZ)$$
$$\geq h(XYZ) + h(Y) + h(XZ)$$

## A Shannon Inequality

Example

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

$$\underline{h(XY) + h(YZ)} + h(XZ)$$
$$\geq h(XYZ) + \underline{h(Y) + h(XZ)}$$

## A Shannon Inequality

**Example**

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

$$\underline{h(XY) + h(YZ)} + h(XZ)$$
$$\geq h(XYZ) + \underline{h(Y) + h(XZ)}$$
$$\geq 2h(XYZ) + h(\emptyset)$$

## A Shannon Inequality

### Example

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

$$\underline{h(XY) + h(YZ)} + h(XZ)$$
$$\geq h(XYZ) + \underline{h(Y) + h(XZ)}$$
$$\geq 2h(XYZ) + h(\emptyset)$$
$$= 2h(XYZ)$$

## A Shannon Inequality

### Example

$h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$

$$\underline{h(XY) + h(YZ)} + h(XZ)$$
$$\geq h(XYZ) + \underline{h(Y) + h(XZ)}$$
$$\geq 2h(XYZ) + h(\emptyset)$$
$$= 2h(XYZ)$$

Note: $X$ is covered 2 times in each expressions. Same for $Y$, same for $Z$.

Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

From Query to Information Inequality:

#### Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X), \qquad |Q| \leq (|R| \cdot |S| \cdot |T|)^{1/2}.$

## Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

From Query to Information Inequality:

### Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X), \qquad |Q| \leq (|R| \cdot |S| \cdot |T|)^{1/2}.$

Instance $\boldsymbol{D} = (R^D, S^D, T^D); \quad p : Q(\boldsymbol{D}) \rightarrow [0, 1]$ uniform; $\boldsymbol{h}$ its entropy.

## Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \le |R_1|^{w_1} \cdots |R_m|^{w_m}}$

From Query to Information Inequality:

### Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X), \qquad |Q| \le (|R| \cdot |S| \cdot |T|)^{1/2}.$

Instance $\boldsymbol{D} = (R^D, S^D, T^D); \qquad p : Q(\boldsymbol{D}) \to [0, 1]$ uniform; $\boldsymbol{h}$ its entropy.

$$\log |R^D| + \log |S^D| + \log |T^D|$$
$$\ge h(XY) + h(YZ) + h(XZ)$$

## Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

From Query to Information Inequality:

### Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X), \qquad |Q| \leq (|R| \cdot |S| \cdot |T|)^{1/2}.$

Instance $\boldsymbol{D} = (R^D, S^D, T^D)$; $\quad p : Q(\boldsymbol{D}) \to [0, 1]$ uniform; $\quad \boldsymbol{h}$ its entropy.

$$\log |R^D| + \log |S^D| + \log |T^D|$$
$$\geq h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$$

# Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

From Query to Information Inequality:

## Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X), \qquad |Q| \leq (|R| \cdot |S| \cdot |T|)^{1/2}.$

Instance $\boldsymbol{D} = (R^D, S^D, T^D);$ $\quad p : Q(\boldsymbol{D}) \rightarrow [0, 1]$ uniform; $\quad \boldsymbol{h}$ its entropy.

$$\begin{aligned} \log |R^D| + \log |S^D| + \log |T^D| \\ \geq h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ) \\ = 2\log |Q(\boldsymbol{D})| \end{aligned}$$

## Proof of the AGM Upper Bound: Part 1: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

From Query to Information Inequality:

### Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X), \qquad |Q| \leq (|R| \cdot |S| \cdot |T|)^{1/2}.$

Instance $\boldsymbol{D} = (R^D, S^D, T^D); \quad p : Q(\boldsymbol{D}) \to [0, 1]$ uniform; $\boldsymbol{h}$ its entropy.

$$\log |R^D| + \log |S^D| + \log |T^D|$$
$$\geq h(XY) + h(YZ) + h(XZ) \geq 2h(XYZ)$$
$$= 2 \log |Q(\boldsymbol{D})|$$

For a general query $Q(\boldsymbol{X}) = R_1(\boldsymbol{Y}_1) \wedge \cdots \wedge R_m(\boldsymbol{Y}_m)$:

If $\boxed{\sum_j w_j h(\boldsymbol{Y}_j) \geq h(\boldsymbol{X})}$ then $\boxed{|R_1|^{w_1} \cdots |R_m|^{w_m} \geq |Q|}$

## Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

Theorem (Generalized Shearer's Inequality)

If $\boldsymbol{w}$ is a frac. edge cover, then $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

# Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

Theorem (Generalized Shearer's Inequality)

*If $\boldsymbol{w}$ is a frac. edge cover, then* $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:

$$\boxed{k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})}$$

Each variable is "covered $\geq k_0$ times".

## Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

Theorem (Generalized Shearer's Inequality)

If $\boldsymbol{w}$ is a frac. edge cover, then $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:
$$\boxed{k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})}$$
Each variable is "covered $\geq k_0$ times".

Repeatedly rewrite $h(\boldsymbol{Y}_i) + h(\boldsymbol{Y}_j) \to h(\boldsymbol{Y}_i \cup \boldsymbol{Y}_j) + h(\boldsymbol{Y}_i \cap \boldsymbol{Y}_j)$

## Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

**Theorem (Generalized Shearer's Inequality)**

If $\boldsymbol{w}$ is a frac. edge cover, then $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:
$$k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})$$
Each variable is "covered $\geq k_0$ times".

Repeatedly rewrite $h(\boldsymbol{Y}_i) + h(\boldsymbol{Y}_j) \rightarrow h(\boldsymbol{Y}_i \cup \boldsymbol{Y}_j) + h(\boldsymbol{Y}_i \cap \boldsymbol{Y}_j)$

- Every variable remains covered $\geq k_0$ times.

Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

Theorem (Generalized Shearer's Inequality)

If **w** is a frac. edge cover, then $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:
$$\boxed{k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})}$$
Each variable is "covered $\geq k_0$ times".

Repeatedly rewrite $h(\boldsymbol{Y}_i) + h(\boldsymbol{Y}_j) \rightarrow h(\boldsymbol{Y}_i \cup \boldsymbol{Y}_j) + h(\boldsymbol{Y}_i \cap \boldsymbol{Y}_j)$

- Every variable remains covered $\geq k_0$ times.
- $\sum_\ell |\boldsymbol{Y}_\ell|^2$ strictly increases (homework!).

## Proof of the AGM Upper Bound: Part 2:   $|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}$

### Theorem (Generalized Shearer's Inequality)

*If $\boldsymbol{w}$ is a frac. edge cover, then* $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:
$$k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})$$
Each variable is "covered $\geq k_0$ times".

Repeatedly rewrite $h(\boldsymbol{Y}_i) + h(\boldsymbol{Y}_j) \rightarrow h(\boldsymbol{Y}_i \cup \boldsymbol{Y}_j) + h(\boldsymbol{Y}_i \cap \boldsymbol{Y}_j)$

- Every variable remains covered $\geq k_0$ times.
- $\sum_\ell |\boldsymbol{Y}_\ell|^2$ strictly increases (homework!).      When do we stop?

# Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

Theorem (Generalized Shearer's Inequality)

*If $\boldsymbol{w}$ is a frac. edge cover, then* $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:

$$k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})$$

Each variable is "covered $\geq k_0$ times".

Repeatedly rewrite $h(\boldsymbol{Y}_i) + h(\boldsymbol{Y}_j) \to h(\boldsymbol{Y}_i \cup \boldsymbol{Y}_j) + h(\boldsymbol{Y}_i \cap \boldsymbol{Y}_j)$

- Every variable remains covered $\geq k_0$ times.

- $\sum_\ell |\boldsymbol{Y}_\ell|^2$ strictly increases (homework!).              When do we stop?

- We stop when $\boldsymbol{Y}_1 \supset \boldsymbol{Y}_2 \supset \cdots$

# Proof of the AGM Upper Bound: Part 2: $\boxed{|Q| \leq |R_1|^{w_1} \cdots |R_m|^{w_m}}$

## Theorem (Generalized Shearer's Inequality)

*If $\boldsymbol{w}$ is a frac. edge cover, then* $\boxed{w_1 h(\boldsymbol{Y}_1) + \cdots + w_m h(\boldsymbol{Y}_m) \geq h(\boldsymbol{X})}$.

**Proof** following [Balister and Bollobás, 2012]. Convert to integers:

$$k_1 h(\boldsymbol{Y}_1) + \cdots + k_m h(\boldsymbol{Y}_m) \geq k_0 h(\boldsymbol{X})$$

Each variable is "covered $\geq k_0$ times".

Repeatedly rewrite $h(\boldsymbol{Y}_i) + h(\boldsymbol{Y}_j) \rightarrow h(\boldsymbol{Y}_i \cup \boldsymbol{Y}_j) + h(\boldsymbol{Y}_i \cap \boldsymbol{Y}_j)$

- Every variable remains covered $\geq k_0$ times.

- $\sum_\ell |\boldsymbol{Y}_\ell|^2$ strictly increases (homework!).      When do we stop?

- We stop when $\boldsymbol{Y}_1 \supset \boldsymbol{Y}_2 \supset \cdots$

- Thus, $\boldsymbol{Y}_1 = \boldsymbol{X}$ and $k_1 \geq k_0$.      $\boxed{\cdots \geq k_1 h(\boldsymbol{Y}_1) \geq k_0 h(\boldsymbol{X})}$

This completes the proof of the Upper AGM Bound.

## Discussion

- Shearer's inequality: apply submodularity repeatedly, in any order!

- Shearer inequalities correspond 1-1 to fractional edge covers.

- Any inequality is an upper bound on $|Q|$: $AGM(Q)$ is the smallest.

- How tight is $AGM(Q)$ upper bound? Next

# Proof of the Lower Bound

## Proof of the AGM Lower Bound

By example:
$$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X) \qquad AGM(Q) = \min_{\boldsymbol{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$$

## Proof of the AGM Lower Bound

By example:
$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$     $AGM(Q) = \min_{\boldsymbol{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$

---

**Primal program:**
Minimize
$w_R \log |R| + w_S \log |S| + w_T \log |T|$
where $\boldsymbol{w}$ is frac. edge cover:

| | | | | |
|---|---|---|---|---|
| $X :$ | $w_R +$ | | $w_T$ | $\geq 1$ |
| $Y :$ | $w_R +$ | $w_S$ | | $\geq 1$ |
| $Z :$ | | $w_S +$ | $w_T$ | $\geq 1$ |

---

## Proof of the AGM Lower Bound

By example:
$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$          $AGM(Q) = \min_{\mathbf{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$

**Primal program:**
Minimize
$w_R \log |R| + w_S \log |S| + w_T \log |T|$
where $\mathbf{w}$ is frac. edge cover:

| $X:$ | $w_R+$ | | $w_T$ | $\geq 1$ |
|---|---|---|---|---|
| $Y:$ | $w_R+$ | $w_S$ | | $\geq 1$ |
| $Z:$ | | $w_S+$ | $w_T$ | $\geq 1$ |

**Dual program:**
Maximize
$v_X + v_Y + v_Z$
where $\mathbf{v}$ is "frac. vertex packing":

| $R:$ | $v_X+$ | $v_Y$ | | $\leq \log |R|$ |
|---|---|---|---|---|
| $S:$ | | $v_Y+$ | $v_Z$ | $\leq \log |S|$ |
| $T:$ | $v_X+$ | | $v_Z$ | $\leq \log |T|$ |

## Proof of the AGM Lower Bound

By example:
$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$       $AGM(Q) = \min_{\boldsymbol{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$

**Primal program:**
Minimize
$w_R \log |R| + w_S \log |S| + w_T \log |T|$
where $\boldsymbol{w}$ is frac. edge cover:

| | | | | |
|---|---|---|---|---|
| $X$ : | $w_R+$ | | $w_T$ | $\geq 1$ |
| $Y$ : | $w_R+$ | $w_S$ | | $\geq 1$ |
| $Z$ : | | $w_S+$ | $w_T$ | $\geq 1$ |

**Dual program:**
Maximize
$v_X + v_Y + v_Z$
where $\boldsymbol{v}$ is "frac. vertex packing":

| | | | | |
|---|---|---|---|---|
| $R$ : | $v_X+$ | $v_Y$ | | $\leq \log |R|$ |
| $S$ : | | $v_Y+$ | $v_Z$ | $\leq \log |S|$ |
| $T$ : | $v_X+$ | | $v_Z$ | $\leq \log |T|$ |

Take optimum $\boldsymbol{v}$, define: $\text{Dom}(X) \stackrel{\text{def}}{=} [\lfloor 2^{v_X} \rfloor]$, $\text{Dom}(Y) \stackrel{\text{def}}{=} [\lfloor 2^{v_Y} \rfloor]$, $\text{Dom}(Z) \stackrel{\text{def}}{=} [\lfloor 2^{v_Z} \rfloor]$.

## Proof of the AGM Lower Bound

By example:
$$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X) \qquad AGM(Q) = \min_{\boldsymbol{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$$

<div>

**Primal program:**
Minimize
$w_R \log |R| + w_S \log |S| + w_T \log |T|$
where $\boldsymbol{w}$ is frac. edge cover:

$$
\begin{array}{ccccc}
X : & w_R + & & w_T & \geq 1 \\
Y : & w_R + & w_S & & \geq 1 \\
Z : & & w_S + & w_T & \geq 1
\end{array}
$$

</div>

<div>

**Dual program:**
Maximize
$$v_X + v_Y + v_Z$$
where $\boldsymbol{v}$ is "frac. vertex packing":

$$
\begin{array}{ccccc}
R : & v_X + & v_Y & & \leq \log |R| \\
S : & & v_Y + & v_Z & \leq \log |S| \\
T : & v_X + & & v_Z & \leq \log |T|
\end{array}
$$

</div>

Take optimum $\boldsymbol{v}$, define: $\text{Dom}(X) \stackrel{\text{def}}{=} [\lfloor 2^{v_X} \rfloor]$, $\text{Dom}(Y) \stackrel{\text{def}}{=} [\lfloor 2^{v_Y} \rfloor]$, $\text{Dom}(Z) \stackrel{\text{def}}{=} [\lfloor 2^{v_Z} \rfloor]$.

Worst-case instance (cartesian products): $R^* \stackrel{\text{def}}{=} \text{Dom}(X) \times \text{Dom}(Y)$, $S^*, T^* \stackrel{\text{def}}{=} \cdots$

## Proof of the AGM Lower Bound

By example:
$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$     $AGM(Q) = \min_{\boldsymbol{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$

---

**Primal program:**
Minimize
$w_R \log |R| + w_S \log |S| + w_T \log |T|$
where $\boldsymbol{w}$ is frac. edge cover:

$$
\begin{array}{lllll}
X: & w_R+ & & w_T & \geq 1 \\
Y: & w_R+ & w_S & & \geq 1 \\
Z: & & w_S+ & w_T & \geq 1
\end{array}
$$

---

**Dual program:**
Maximize
$$v_X + v_Y + v_Z$$
where $\boldsymbol{v}$ is "frac. vertex packing":

$$
\begin{array}{lllll}
R: & v_X+ & v_Y & & \leq \log |R| \\
S: & & v_Y+ & v_Z & \leq \log |S| \\
T: & v_X+ & & v_Z & \leq \log |T|
\end{array}
$$

---

Take optimum $\boldsymbol{v}$, define: $\text{Dom}(X) \stackrel{\text{def}}{=} [\lfloor 2^{v_X} \rfloor]$, $\text{Dom}(Y) \stackrel{\text{def}}{=} [\lfloor 2^{v_Y} \rfloor]$, $\text{Dom}(Z) \stackrel{\text{def}}{=} [\lfloor 2^{v_Z} \rfloor]$.

Worst-case instance (cartesian products): $R^* \stackrel{\text{def}}{=} \text{Dom}(X) \times \text{Dom}(Y)$, $S^*, T^* \stackrel{\text{def}}{=} \cdots$

$|Q^*| = \lfloor 2^{v_X} \rfloor \cdot \lfloor 2^{v_Y} \rfloor \cdot \lfloor 2^{v_Z} \rfloor \geq \frac{1}{8} 2^{v_X + v_Y + v_Z}$

## Proof of the AGM Lower Bound

By example:
$R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$     $AGM(Q) = \min_{\mathbf{w}} |R|^{w_R} \cdot |S|^{w_S} \cdot |T|^{w_T}$

---

**Primal program:**
Minimize
$w_R \log |R| + w_S \log |S| + w_T \log |T|$
where $\mathbf{w}$ is frac. edge cover:

| | | | | |
|---|---|---|---|---|
| $X:$ | $w_R+$ | | $w_T$ | $\geq 1$ |
| $Y:$ | $w_R+$ | $w_S$ | | $\geq 1$ |
| $Z:$ | | $w_S+$ | $w_T$ | $\geq 1$ |

---

**Dual program:**
Maximize
$$v_X + v_Y + v_Z$$
where $\mathbf{v}$ is "frac. vertex packing":

| | | | | |
|---|---|---|---|---|
| $R:$ | $v_X+$ | $v_Y$ | | $\leq \log |R|$ |
| $S:$ | | $v_Y+$ | $v_Z$ | $\leq \log |S|$ |
| $T:$ | $v_X+$ | | $v_Z$ | $\leq \log |T|$ |

---

Take optimum $\mathbf{v}$, define: $\text{Dom}(X) \stackrel{\text{def}}{=} [\lfloor 2^{v_X} \rfloor], \text{Dom}(Y) \stackrel{\text{def}}{=} [\lfloor 2^{v_Y} \rfloor], \text{Dom}(Z) \stackrel{\text{def}}{=} [\lfloor 2^{v_Z} \rfloor].$

Worst-case instance (cartesian products): $R^* \stackrel{\text{def}}{=} \text{Dom}(X) \times \text{Dom}(Y), S^*, T^* \stackrel{\text{def}}{=} \cdots$

$|Q^*| = \lfloor 2^{v_X} \rfloor \cdot \lfloor 2^{v_Y} \rfloor \cdot \lfloor 2^{v_Z} \rfloor \geq \frac{1}{8} 2^{v_X + v_Y + v_Z} = \frac{1}{8} 2^{w_1^* \log |R| + w_2^* \log |S| + w_3^* \log |T|} = \frac{1}{8} 2^{AGM(Q)}$

# Special Case: $|R| = |S| = \cdots = N$

### Definition

Fix a hypergraph $(V, E)$; $(v_X)_{X \in V} \in \mathbb{R}_+^{|V|}$ is a fractional vertex packing if:
$$\forall \boldsymbol{Y} \in E :, \quad \boxed{\sum_{X \in \boldsymbol{Y}} v_X \leq 1}$$

# Special Case: $|R| = |S| = \cdots = N$

## Definition

Fix a hypergraph $(V, E)$; $(v_X)_{X \in V} \in \mathbb{R}_+^{|V|}$ is a fractional vertex packing if:
$$\forall \mathbf{Y} \in E :, \boxed{\sum_{X \in \mathbf{Y}} v_X \leq 1}$$

When $|R| = |S| = \cdots = N$, then replace

$$\begin{array}{c} v_R + v_S \leq \log N \\ v_R + v_T \leq \log N \\ \cdots \end{array}$$

with

$$\begin{array}{c} v_R + v_S \leq 1 \\ v_R + v_T \leq 1 \\ \cdots \end{array}$$

times $\log N$.

# Special Case: $|R| = |S| = \cdots = N$

> **Definition**
>
> Fix a hypergraph $(V, E)$; $(v_X)_{X \in V} \in \mathbb{R}_+^{|V|}$ is a fractional vertex packing if:
> $$\forall \boldsymbol{Y} \in E :, \quad \boxed{\sum_{X \in \boldsymbol{Y}} v_X \leq 1}$$

When $|R| = |S| = \cdots = N$, then replace

$$
\begin{array}{c}
v_R + v_S \leq \log N \\
v_R + v_T \leq \log N \\
\cdots
\end{array}
$$

with

$$
\begin{array}{c}
v_R + v_S \leq 1 \\
v_R + v_T \leq 1 \\
\cdots
\end{array}
$$

times $\log N$.

Then: $\qquad R = [N^{v_X}] \times [N^{v_Y}], \ S = [N^{v_Y}] \times [N^{v_Z}], \ T = [N^{v_X}] \times [N^{v_Z}].$

$$Q = [N^{v_X}] \times [N^{v_Y}] \times [N^{v_Z}]$$

Background
○○○○
Output Bound
○○○
AGM Bound
○○○○○
Proof: Upper Bound
○○○○○○○
Proof: Lower Bound
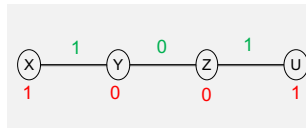○○○●○
Extensions
○○○○○○○○○○○

# Examples

$|R| = |S| = \cdots = N$

$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$

## Examples

$|R| = |S| = \cdots = N$

$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$

## Examples

$|R| = |S| = \cdots = N$

$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$
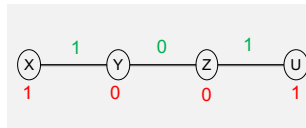$R = [N] \times [1], \ S = [1] \times [1], \ T = [1] \times [N].$

# Examples

$|R| = |S| = \cdots = N$

$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$
$R = [N] \times [1], \; S = [1] \times [1], \; T = [1] \times [N].$



$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge K(U, V)$

## Examples

$|R| = |S| = \cdots = N$

$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$
$R = [N] \times [1], \; S = [1] \times [1], \; T = [1] \times [N].$



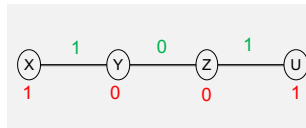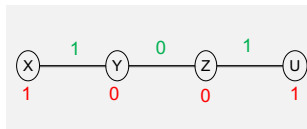$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge K(U, V)$

# Examples

$|R| = |S| = \cdots = N$

$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U)$
$R = [N] \times [1],\ S = [1] \times [1],\ T = [1] \times [N].$



$R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge K(U, V)$
$R = T = [N] \times [1],\ S = K = [1] \times [N]$

## Summary of the AGM Bound

- Upper / lower bound: fractional  edge cover / vertex packing.

- Their equality follows from strong duality.

- The worst-case instance of the AGM bound is a Product Database.

- Full CQs only. Otherwise, ignore non-head variables.

Limitation of AGM: only cardinalities. Next: extensions to other stats.

# Extensions of the AGM Bound

# Simple Functional Dependencies

Given functional dependencies, query output is $\ll$ AGM bound.

Example: $R(X, Y) \wedge S(Y, Z)$: $N^2$ becomes $N$ when $Y \to Z$.

An FD $\boldsymbol{U} \to \boldsymbol{V}$ is simple if $\boldsymbol{U}$ is a single variable.

Method [Khamis et al., 2016]:

- Expand $Q$ to $Q^+$ by replacing each atom $R(\boldsymbol{Y})$ with $R'(\boldsymbol{Y^+})$.

- Compute the AGM bound of $Q^+$.

- This bound is tight. Proof: very useful exercise.

## Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$
Fractional edge covers: $(1, 1, 0), (1, 0, 1), (0, 1, 1), (1/2, 1/2, 1/2)$

$$|Q| \leq \min(|R| \cdot |S|, |R| \cdot |T|, |S| \cdot |T|, \sqrt{|R| \cdot |S| \cdot |T|})$$

## Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$
Fractional edge covers: $(1, 1, 0), (1, 0, 1), (0, 1, 1), (1/2, 1/2, 1/2)$

$$\boxed{|Q| \leq \min(|R| \cdot |S|, |R| \cdot |T|, |S| \cdot |T|, \sqrt{|R| \cdot |S| \cdot |T|})}$$

Assume that $S.Y$ is a key:           $Y \rightarrow Z$

## Example

$Q(X, Y, Z) = R(X, Y) \wedge S(Y, Z) \wedge T(Z, X)$
Fractional edge covers: $(1, 1, 0), (1, 0, 1), (0, 1, 1), (1/2, 1/2, 1/2)$

$$|Q| \leq \min(|R| \cdot |S|, |R| \cdot |T|, |S| \cdot |T|, \sqrt{|R| \cdot |S| \cdot |T|})$$

Assume that $S.Y$ is a key: $\qquad\qquad Y \to Z$
$Q^+(X, Y, Z) = R'(X, Y, Z) \wedge S(Y, Z) \wedge T(Z, X)$
Fractional edge covers: $(1, 0, 0), (0, 1, 1)$

$$|Q| \leq \min(|R|, |S| \cdot |T|)$$

## Discussion

The expansion procedure is very easy, but limited only to simple FDs:

$AGM(Q^+)$ is always an upper bound on $Q$'s output, but may not be tight.

Example

$$Q(X, Y, Z, U) = R(X, Y) \land S(Y, Z) \land T(Z, U) \land A(X, Z, U) \land B(X, Y, U)$$

$$A: XZ \to U; \qquad\qquad B: YU \to X$$

Expansion is useless ($Q^+ = Q$).

## More Statistics

Statistics for a relation $R(U, V, W, \ldots)$:

- Its cardinality $|R|$.

- Number distinct values of an attribute / set of attributes, e.g. $|R.X|$.

- Max degree of an attribute / set of attributes, e.g. $\max(\deg_R(VW|U))$

- The max degree of a projection, e.g. $\max(\deg_R(V|U))$.

- The $\ell_p$-norm of some degree sequence, e.g. $||\deg_R(V|U)||_2$.

Will use entropic inequalities, beyond Shearer

## Example

$$R = \begin{array}{|ccc|} \hline U & V & W \\ \hline a & 1 & m \\ a & 1 & n \\ a & 2 & m \\ a & 3 & m \\ b & 1 & m \\ b & 5 & m \\ \hline \end{array}$$

$$|R| = 6$$
$$|R.U| = 2$$
$$|R.V| = 4$$
$$|R.UV| = 5$$

$\max(\deg_R(VW|U)) = 4$

$\max(\deg_R(V|U)) = 3$

# Conditional Entropy

The *Conditional Entropy*

$$h(\boldsymbol{V}|\boldsymbol{U}) \stackrel{\text{def}}{=} h(\boldsymbol{U}\boldsymbol{V}) - h(\boldsymbol{U})$$

## Conditional Entropy

The *Conditional Entropy*

$$h(\boldsymbol{V}|\boldsymbol{U}) \stackrel{\text{def}}{=} h(\boldsymbol{U}\boldsymbol{V}) - h(\boldsymbol{U})$$

What it means: $h(\boldsymbol{V}|\boldsymbol{U}) = \mathbb{E}_{\boldsymbol{u}}[h(\boldsymbol{V}|\boldsymbol{U} = \boldsymbol{u})]$

## Conditional Entropy

The *Conditional Entropy*

$$h(\boldsymbol{V}|\boldsymbol{U}) \stackrel{\text{def}}{=} h(\boldsymbol{UV}) - h(\boldsymbol{U})$$

What it means: $h(\boldsymbol{V}|\boldsymbol{U}) = \mathbb{E}_{\boldsymbol{u}}[h(\boldsymbol{V}|\boldsymbol{U} = \boldsymbol{u})]$

The submodularity inequality can be written equivalently as:

$$h(\boldsymbol{V}|\boldsymbol{U}) \geq h(\boldsymbol{V}|\boldsymbol{UW})$$

## From Entropy to Statistics

Fix a joint probability distribution of the variables $\boldsymbol{X}$, with support $R(\boldsymbol{X})$:

$$h(\boldsymbol{X}) \leq \log |R|$$

$$h(\boldsymbol{V}|\boldsymbol{U}) \leq \log \left(\max \deg_R(\boldsymbol{V}|\boldsymbol{U})\right)$$

$$h(\boldsymbol{UV}) + (p-1)h(\boldsymbol{V}|\boldsymbol{V}) \leq \log ||\deg_R(\boldsymbol{V}|\boldsymbol{U})||_p^p$$

## Example

$$Q = R(X, Y) \land S(Y, Z) \land T(Z, U) \land A(X, Z, U) \land B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$:  $\hspace{6cm}$ $AGM(Q) = N^2$.

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$: $\qquad\qquad\qquad\qquad\qquad AGM(Q) = N^2.$

If the FDs $XZ \to U$ and $YU \to X$ hold: $\qquad\qquad\qquad\qquad |Q| \leq N^{3/2}.$

$\log |R| + \log |S| + \log |T| + \log \max \deg_A(U|XZ) + \log \max \deg_B(X|YU) \geq$

$\geq h(XY) + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$: $\hspace{4cm} AGM(Q) = N^2.$
If the FDs $XZ \rightarrow U$ and $YU \rightarrow X$ hold: $\hspace{2.5cm} |Q| \leq N^{3/2}.$

$\log|R| + \log|S| + \log|T| + \log \max \deg_A(U|XZ) + \log \max \deg_B(X|YU) \geq$
$\geq \underline{h(XY)} + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$: $\qquad\qquad\qquad\qquad\qquad AGM(Q) = N^2$.

If the FDs $XZ \to U$ and $YU \to X$ hold: $\qquad\qquad\qquad\qquad |Q| \leq N^{3/2}$.

$\log|R| + \log|S| + \log|T| + \log\max\deg_A(U|XZ) + \log\max\deg_B(X|YU) \geq$

$\geq \underline{h(XY) + h(YZ)} + h(ZU) + h(U|XZ) + h(X|YU)$

$\qquad \geq h(XYZ) + h(Y) + h(ZU) + h(U|XZ) + h(X|YU)$

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$: $\hspace{4cm} AGM(Q) = N^2$.
If the FDs $XZ \rightarrow U$ and $YU \rightarrow X$ hold: $\hspace{2.5cm} |Q| \leq N^{3/2}$.

$\log |R| + \log |S| + \log |T| + \log \max \deg_A(U|XZ) + \log \max \deg_B(X|YU) \geq$
$\geq h(XY) + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$
$\hspace{1cm} \geq h(XYZ) + \underline{h(Y) + h(ZU)} + h(U|XZ) + h(X|YU)$

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$:  $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad AGM(Q) = N^2$.
If the FDs $XZ \rightarrow U$ and $YU \rightarrow X$ hold:  $\quad\quad\quad\quad\quad\quad\quad\quad |Q| \leq N^{3/2}$.

$\log|R| + \log|S| + \log|T| + \log \max \deg_A(U|XZ) + \log \max \deg_B(X|YU) \geq$
$\geq h(XY) + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$
$\quad\quad \geq h(XYZ) + \underline{h(Y) + h(ZU)} + h(U|XZ) + h(X|YU)$
$\quad\quad \geq h(XYZ) + h(YZU) + h(U|XZ) + h(X|YU)$

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$:  $\qquad\qquad\qquad\qquad\qquad AGM(Q) = N^2$.
If the FDs $XZ \rightarrow U$ and $YU \rightarrow X$ hold: $\qquad\qquad\qquad |Q| \leq N^{3/2}$.

$\log|R| + \log|S| + \log|T| + \log\max\deg_A(U|XZ) + \log\max\deg_B(X|YU) \geq$
$\geq h(XY) + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$
$\qquad \geq h(XYZ) + h(Y) + h(ZU) + h(U|XZ) + h(X|YU)$
$\qquad \geq h(XYZ) + h(YZU) + \underline{h(U|XZ)} + \underline{h(X|YU)}$

## Example

$$Q = R(X, Y) \wedge S(Y, Z) \wedge T(Z, U) \wedge A(X, Z, U) \wedge B(X, Y, U)$$

Assume $|R| = |S| = |T| = N$: $\qquad AGM(Q) = N^2$.

If the FDs $XZ \to U$ and $YU \to X$ hold: $\qquad |Q| \leq N^{3/2}$.

$\log |R| + \log |S| + \log |T| + \log \max \deg_A(U|XZ) + \log \max \deg_B(X|YU) \geq$

$\geq h(XY) + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$

$\qquad \geq h(XYZ) + h(Y) + h(ZU) + h(U|XZ) + h(X|YU)$

$\qquad \geq h(XYZ) + h(YZU) + \underline{h(U|XZ)} + \underline{h(X|YU)}$

$\qquad \geq h(XYZ) + h(YZU) + h(U|XYZ) + h(X|YZU)$

$\qquad = 2h(XYZU) = 2 \log |Q|$

## Example

$$Q = R(X,Y) \wedge S(Y,Z) \wedge T(Z,U) \wedge A(X,Z,U) \wedge B(X,Y,U)$$

Assume $|R| = |S| = |T| = N$: $\hspace{4cm} AGM(Q) = N^2$.
If the FDs $XZ \to U$ and $YU \to X$ hold: $\hspace{2.3cm} |Q| \le N^{3/2}$.

$\log |R| + \log |S| + \log |T| + \log \max \deg_A(U|XZ) + \log \max \deg_B(X|YU) \ge$
$\ge h(XY) + h(YZ) + h(ZU) + h(U|XZ) + h(X|YU)$
$\hspace{1.2cm} \ge h(XYZ) + h(Y) + h(ZU) + h(U|XZ) + h(X|YU)$
$\hspace{1.2cm} \ge h(XYZ) + h(YZU) + h(U|XZ) + h(X|YU)$
$\hspace{1.2cm} \ge h(XYZ) + h(YZU) + h(U|XYZ) + h(X|YZU)$
$\hspace{1.2cm} = 2h(XYZU) = 2 \log |Q|$

$$\boxed{|Q| \le \sqrt{|R| \cdot |S| \cdot |T| \cdot \max(\deg(U|XZ)) \cdot \max(\deg(X|YU))}}$$

## Discussion

- AGM/Shearer limited to cardinality statistics.

- More general statistics require general entropic inequalities.

- Everything gets harder: fractional edge cover no longer sufficient, order of the submodularity matters.

- Can we compute the upper bound? Is it tight? Yes and no, it's complicated [Suciu, 2023].

- Do they work in practice? Yes, but you need to do the engineering work [Deeds et al., 2023].

Atserias, A., Grohe, M., and Marx, D. (2013).
Size bounds and query plans for relational joins.
*SIAM J. Comput.*, 42(4):1737–1767.

Balister, P. and Bollobás, B. (2012).
Projections, entropy and sumsets.
*Comb.*, 32(2):125–141.

Deeds, K. B., Suciu, D., and Balazinska, M. (2023).
Safebound: A practical system for generating cardinality bounds.
*Proc. ACM Manag. Data*, 1(1):53:1–53:26.

Khamis, M. A., Ngo, H. Q., and Suciu, D. (2016).
Computing join queries with functional dependencies.
In Milo, T. and Tan, W., editors, *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 327–342. ACM.

Suciu, D. (2023).
Applications of information inequalities to database theory problems.
In *LICS*, pages 1–30.