# Introduction to Bayesian Analysis with HllSets

From whence do we come, and whither do we tread?
The essence of our being, a riddle left unsaid.
Beneath the endless sky, so many hearts in anguish yearn,
Turning to mere embers, to dust—yet where's the smoke that's spurned? [1]

HllSets [2] seamlessly align with the Bayesian framework of probabilistic relationships between events or entities, which can be represented as sets. In this context, we will define events as the occurrence of an entity instance at a specific point in time. This approach enables us to use the terms "event" and "entity instance" interchangeably.

## SGS universe

The article [3] explores the concept of Self-Generative Systems (SGS), inspired by John von Neumann's theory of self-reproducing automata. It applies this concept to a Metadata Management System that leverages HyperLogLog Sets (HllSets) [2] and graph databases to efficiently manage data and its associated metadata.

SGS content is inherently linked to the ingested data and is manifested as a collection of entities. Each entity is essentially a shell encapsulating HllSets that comprehensively define the entity instance.

We define the SGS universe as the aggregate of all current entity instances. Each SGS has a unique universe, which, unlike a global universe, is subject to change. We address these changes in the same manner as we handle modifications to any other entity instances.

In [3], we identified a self-reproducing loop triggered by executing the Commit command. To recap, the Commit command transfers ingested data, already transformed into entity instances, from the staging area to the forefront of the SGS.

There are two key takeaways from this process:
1. By the time we are prepared to commit newly acquired data, we possess a collection of entity instances along with its local universe—a union of all entities in the staging area.
2. The Commit action does not overwrite the existing SGS head or the current SGS universe. Instead, if new data introduces changes to the head data, the Commit command will relocate affected entities to the tail, thereby updating the current head.

**Definition 1**: The current SGS universe is defined as the union of all entities from the HEAD of the SGS:

$$U_{head} = \bigcup_{i \in HEAD} E_i,$$

where: HEAD is a collection of all entities (instances) in the HEAD of the SGS.

**Definition 2**: The SGS universe for the t-slice from the SGS tail is defined as the union of all entities from the t-slice of the SGS tail. The t-slice comprises the entities pushed into the SGS tail by the Commit command executed at time t:

$$U_{t-slice} = \bigcup_{i \in SLICE(t)} E_i,$$

where: SLICE(t) is a collection of all entities (instances) in the SLICE(t) in the tail of the SGS.

We can seamlessly integrate the current SGS universe with sliced universes from the SGS tail, allowing us to extend Bayesian analysis to the history of SGS. Furthermore, our research should not be confined to the present; we should explore the history of SGS, hypothesizing a slice from the historical context relevant to our current research and constructing a Bayesian model from that point.

We emphasize the importance of discussing the concept of the universe because there is no "universal" universe. What was deemed acceptable a few years ago may no longer hold today, and the reasoning of previous generations is not always applicable without scrutiny.

Bayes formulas

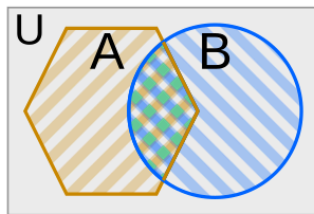The following is a direct copy from [5]:

Bayes' theorem is stated mathematically as the following equation:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

where A and B are events $P(b) \neq 0$.

Below is a graphical presentation from [5].

$$P(A) = \frac{\text{⬡}}{\text{▭}} \quad , \quad P(B|A) = \frac{\text{◈}}{\text{⬡}}$$

$$P(B) = \frac{\text{◯}}{\text{▭}} \quad , \quad P(A|B) = \frac{\text{◈}}{\text{◯}}$$

$$P(A) \cdot P(B|A) = \frac{\text{⬡}}{\text{▭}} \times \frac{\text{◈}}{\text{⬡}} = \frac{\text{◈}}{\text{▭}}$$

$$P(B) \cdot P(A|B) = \frac{\text{◯}}{\text{▭}} \times \frac{\text{◈}}{\text{◯}} = \frac{\text{◈}}{\text{▭}}$$

$$= P(A) \cdot P(B|A) \, , \ \ \text{i.e.}$$



$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$

$$P(B|A) = \frac{P(B) \cdot P(A|B)}{P(A)}$$

In the framework of SGS analysis, it is crucial to substitute any generic gray rectangle with a corresponding "local" universe. For instance, when examining SGS head entities, we should utilize a head-specific universe.

## Example

The following is a demonstration of applying Bayes formulas to HllSets.

```python
Python
include("src/sets.jl")

import .HllSets as set
using Random

A = set.HllSet{10}()
B = set.HllSet{10}()
C = set.HllSet{10}()
```

```
items_t1 = Set(["string0", "string1", "string2", "string3", "string4",
"string5", "string6", "string7", "string8", "string9", "string10"])
items_t2 = Set(["string3", "string4", "string5", "string6", "string7",
"string8", "string9", "string10", "string11"])
items_t3 = Set(["string5", "string6", "string7", "string8", "string9",
"string10", "string11"])

set.add!(A, items_t1)
set.add!(B, items_t2)
set.add!(C, items_t3)

U = A ∪ B ∪ C

println("A: ", count(A))
println("B: ", count(B))
println("C: ", count(C))
println("U: ", count(U), "\n")

println("AB = A ∩ B: ", count(A ∩ B))
println("AC = A ∩ C: ", count(A ∩ C))
println("BC = B ∩ C: ", count(B ∩ C), "\n")

println("P(A | B) = AB / B: ", count(A ∩ B) / count(B))
println("P(B | A) = AB / A: ", count(A ∩ B) / count(A))
println("P(A | C) = AC / C: ", count(A ∩ C) / count(C))
println("P(C | A) = AC / A: ", count(A ∩ C) / count(A), "\n")

println("P(B | C) = BC / C: ", count(B ∩ C) / count(C))
println("P(C | B) = BC / B: ", count(B ∩ C) / count(B), "\n")
```

and here is an output of this code:

```Python
A: 11
B: 9
C: 7
U: 12

AB = A ∩ B: 8
AC = A ∩ C: 6
BC = B ∩ C: 7
```

```
P(A | B) = AB / B: 0.8888888888888888
P(B | A) = AB / A: 0.7272727272727273
P(A | C) = AC / C: 0.8571428571428571
P(C | A) = AC / A: 0.5454545454545454

P(B | C) = BC / C: 1.0
P(C | B) = BC / B: 0.7777777777777778
```

# Summary

This article discusses the integration of HllSets within Bayesian analysis, particularly in the context of Self-Generative Systems (SGS) that manage data using HyperLogLog sets and graph databases. It explores the concept of a "universe" in SGS, where each universe is a collection of entity instances that can change over time. The document also illustrates applying Bayes' theorem using HllSets, demonstrating probabilistic calculations with set operations to manage and analyze large datasets effectively in a dynamic environment.

# References

1. From The Rubaiyat of Omar Khayyam. Translation (with a help from WorkGPT for Workspaces Google) from Russian translation by Osip Rumer (1883-1954). Source: Омар Хайям. Четверостишия. М., 1938 (Omar Khayyám. Rubaiyat. Moscow, 1938).
2. Alex Mylnikov on LinkedIn: HyperLogLog based approximation for very big Datasets
3. HllSet commit and Self Reproductive System | by Alex Mylnikov | Jul, 2024 | Medium
4. How To Reason Mathematically
5. Bayes' theorem - Wikipedia