# Generalization: Moments

- **Suppose a stream has elements chosen from a set _A_ of _N_ values**

- **Let $m_i$ be the number of times value _i_ occurs in the stream**

- **The $k^{\text{th}}$ _moment_ is**

$$\sum_{i \in A} (m_i)^k$$

# Special Cases

$$\sum_{i \in A} (m_i)^k$$

- **0th moment =** number of distinct elements
  - The problem just considered
- **1st moment =** count of the numbers of elements = length of the stream
  - Easy to compute
- **2nd moment = *surprise number S* =** a measure of how uneven the distribution is

# Example: Surprise Number

- **Stream of length 100**
- **11 distinct values**

- Item counts: **10, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9**
  **Surprise $S$ = 910**

- Item counts: **90, 1, 1, 1, 1, 1, 1, 1 ,1, 1, 1**
  **Surprise $S$ = 8,110**

# AMS Method

- **AMS method works for all moments**
- **Gives an unbiased estimate**
- **We will just concentrate on the 2ⁿᵈ moment $S$**
- **We keep track of many variables $X$:**
  - For each variable $X$ we store $X.el$ and $X.val$
    - $X.el$ corresponds to the item $i$
    - $X.val$ corresponds to the **count** of item $i$
  - Note this requires a count in main memory, so number of $X$s is limited
- **Our goal is to compute $S = \sum_i m_i^2$**

# One Random Variable (X)

- **How to set *X.val* and *X.el*?**

  - Assume stream has length *n* **(we relax this later)**

  - Pick some random time *t* (*t<n*) to start,
    so that any time is equally likely

  - Let at time *t* the stream have item *i*. *We set X.el = i*

  - Then we maintain count *c* (*X.val = c*) of the number
    of *i*s in the stream starting from the chosen time *t*

- **Then the estimate of the 2ⁿᵈ moment ($\sum_i m_i^2$) is:**
$$S = f(X) = n(2 \cdot c - 1)$$

  - Note, we keep track of multiple **X**s, ($X_1$, $X_2$,...$X_k$)
    and our final estimate will be $S = 1/k \sum_j f(X_j)$