

Boolean Tuple Set Algorithm

Ohad Asor

April 16, 2018

We present an algorithm for storing and querying fixed-length vectors of booleans. We use statistical methods in order to prune the search space. We assume the reader is familiar with Least-Squares techniques. We also assume that the queries are fixed, known beforehand, and their result is always desired. Let X be our set of tuples $X = \{\mathbf{x}_i\}_{i=1}^n$ where $X \subset \{\pm 1\}^k$. First we assume that if $\mathbf{t} \in X$ then also $-\mathbf{t} \in X$. This is done by adding 1 as a new element in the beginning of every vector, and makes the distribution of our population to be symmetric, so all odd means (including the mean) are zero. The covariance matrix is therefore:

$$\Sigma_X = \mathbb{E}_{\mathbf{x} \in X} [\mathbf{x}\mathbf{x}^T]$$

Call the quantity $\mathbf{x}^T \Sigma_X^{-1} \mathbf{x}$ the *normalized norm* of \mathbf{x} . Suppose we'd like to find all vectors on our set given some of the bits but not all, e.g. all vectors with first and third bit set to one. Denote this query by \mathbf{q} such that $q_i = \pm 1$ if we'd like that bit to be set/unset, or 0 if we don't care. Then

$$\left\| \sqrt{\Sigma_X^{-1}} \mathbf{q} \right\| \tag{1}$$

is the minimum distance of any matching vector from the mean. By that we pruned the search space if we store the vectors sorted by to their norm, as a function of the number of bits given, since the more given bits, the larger is (1). Moreover, the larger the query's norm is, the even larger the contribution to the pruning space, as the size of the sets sorted by distance goes down linearly with the norm, as Chebyshev's inequality would suggest. Observe that we couldn't possibly find a better distance function from the point of view of the first three moments, not only two.

For each new batch of bitvectors to add to X , we need to calculate X' and store its vectors sorted according to their norm, as well as resolve the queries (which we assumed fixed). This is done in a single separate pass over X . We then renormalize the queries, and with a second pass over the data we check only the tuples with high enough norm, as the queries dictate.

We're left to point how to compute the normalized norms online. Sherman-Morrison update is

$$(\Sigma + \mathbf{x}\mathbf{x}^T)^{-1} = \Sigma^{-1} - \frac{\Sigma^{-1} \mathbf{x}\mathbf{x}^T \Sigma^{-1}}{1 + \mathbf{x}^T \Sigma^{-1} \mathbf{x}}$$

so the norm update is:

$$\mathbf{t}^T (\Sigma + \mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} = \mathbf{t}^T \Sigma^{-1} \mathbf{t} - \frac{\mathbf{t}^T \Sigma^{-1} \mathbf{x} \mathbf{x}^T \Sigma^{-1} \mathbf{t}}{1 + \mathbf{x}^T \Sigma^{-1} \mathbf{x}}$$

where $\mathbf{t}^T \Sigma^{-1} \mathbf{t}$ is just the previous norm and the additive update is

$$\frac{\langle \mathbf{t}^T, \Sigma^{-1} \mathbf{x} \rangle^2}{1 + \|\mathbf{x}\|_{\Sigma^{-1}}^2}$$