

# Downtown Bodega Filter Tuning Guide

Calvin Higgins

February 2025

## 1 Background

Downtown Bodega filters (DBFs) [1] are a secure alternative to Learned Bloom filters (LBFs) [2]. Both data structures require tuning several parameters including a decision threshold  $\tau$  and backup filter false positive rate(s) to attain the desired false positive rate  $\epsilon$ . Parameter tuning also impacts the filter’s memory footprint and, for DBFs, worst-case false positive rate. Tuning procedures for LBFs do not directly translate to DBFs, as they do not balance the tradeoff between expected and worst-case false positive rates. To be practical, DBFs must have a lower expected false positive rate than Naor-Eylon filters [3], an alternative secure Bloom filter variant, and reasonable worst-case false positive rate.

## 2 Problem Statement

Given a desired expected false positive rate  $\epsilon$  and maximum false positive rate  $\epsilon_{\max}$ , we wish to choose the decision threshold  $\tau$ , true positive filter false positive rate  $\text{TP}_{\text{FPR}}$ , and false negative filter false positive rate  $\text{FN}_{\text{FPR}}$  such that the resulting Downtown Bodega filter (DBF) has minimal size in bits  $m$ .

## 3 Results

We give a log-linear time vectorizable algorithm to compute the optimal decision threshold  $\tau$ . First, we compute the optimal backup filter false positive rates, and then optimize the choice of  $\tau$  with dynamic programming.

### 3.1 Backup Filter False Positive Rates

We consider the problem of choosing optimal backup filter false positive rates  $\text{TP}_{\text{FPR}}$  and  $\text{FN}_{\text{FPR}}$ , given a learned model with false positive rate  $M_{\text{FPR}}$ .

From Theorem 8 in [1], we have that

$$\epsilon = M_{\text{FPR}}\text{TP}_{\text{FPR}} + M_{\text{TNR}}\text{FN}_{\text{FPR}} \quad (1)$$

and the total size of the DBF is given by

$$m = -\frac{n}{\ln(2)^2} (M_{\text{TPR}} \ln(\text{TP}_{\text{FPR}}) + M_{\text{FNR}} \ln(\text{FN}_{\text{FPR}})) + M_m \quad (2)$$

where  $n$  is the size of the underlying set and  $M_m$  is the memory consumption of the learned model in bits. Manipulating (1), we find that

$$\text{FN}_{\text{FPR}} = \frac{\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}}{M_{\text{TNR}}} \quad (3)$$

so

$$m = -\frac{n}{\ln(2)^2} \left( M_{\text{TPR}} \ln(\text{TP}_{\text{FPR}}) + M_{\text{FNR}} \ln \left( \frac{\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}}{M_{\text{TNR}}} \right) \right) + M_m \quad (4)$$

by plugging into (2).

We wish to minimize  $m$  by choosing  $\text{TP}_{\text{FPR}}$ , so we differentiate  $m$  with respect to  $\text{TP}_{\text{FPR}}$ .

$$\frac{\partial m}{\partial \text{TP}_{\text{FPR}}} = -\frac{n}{\ln(2)^2} \left( \frac{\partial}{\partial \text{TP}_{\text{FPR}}} M_{\text{TPR}} \ln(\text{TP}_{\text{FPR}}) + \frac{\partial}{\partial \text{TP}_{\text{FPR}}} M_{\text{FNR}} \ln \left( \frac{\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}}{M_{\text{TNR}}} \right) \right) \quad (5)$$

$$= -\frac{n}{\ln(2)^2} \left( \frac{M_{\text{TPR}}}{\text{TP}_{\text{FPR}}} + \frac{M_{\text{FNR}}}{\frac{\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}}{M_{\text{TNR}}}} \cdot \frac{M_{\text{FPR}}}{M_{\text{TNR}}} \right) \quad (6)$$

$$= -\frac{n}{\ln(2)^2} \left( \frac{M_{\text{TPR}}}{\text{TP}_{\text{FPR}}} + \frac{M_{\text{FNR}} M_{\text{FPR}}}{\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}} \right) \quad (7)$$

$$(8)$$

Then we set to zero and solve

$$-\frac{n}{\ln(2)^2} \left( \frac{M_{\text{TPR}}}{\text{TP}_{\text{FPR}}} + \frac{M_{\text{FNR}} M_{\text{FPR}}}{\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}} \right) = 0 \quad (9)$$

$$-M_{\text{TPR}} (\epsilon - M_{\text{FPR}} \text{TP}_{\text{FPR}}) - \text{TP}_{\text{FPR}} M_{\text{FNR}} M_{\text{FPR}} = 0 \quad (10)$$

$$(M_{\text{TPR}} - M_{\text{FNR}}) M_{\text{FPR}} \text{TP}_{\text{FPR}} = \epsilon M_{\text{TPR}} \quad (11)$$

$$M_{\text{FPR}} \text{TP}_{\text{FPR}} = \epsilon M_{\text{TPR}} \quad (12)$$

$$\text{TP}_{\text{FPR}} = \epsilon \frac{M_{\text{TPR}}}{M_{\text{FPR}}} \quad (13)$$

Substituting into (3), we find that

$$\text{FN}_{\text{FPR}} = \epsilon \frac{M_{\text{FNR}}}{M_{\text{TNR}}} \quad (14)$$

### 3.2 Learned Model False Positive Rate

In Section 3.1, we found the optimal backup filter false positive rates given a fixed decision threshold  $\tau$ . Now we consider the problem of choosing the optimal decision threshold.

Plugging (13) and (14) into (2), we find that

$$m = -\frac{n}{\ln(2)^2} \left( M_{\text{TPR}} \ln \left( \epsilon \frac{M_{\text{TPR}}}{M_{\text{FPR}}} \right) + M_{\text{FNR}} \ln \left( \epsilon \frac{M_{\text{FNR}}}{M_{\text{TNR}}} \right) \right) + M_m \quad (15)$$

$$= -\frac{n}{\ln(2)^2} \left( M_{\text{TPR}} \ln \left( \frac{M_{\text{TPR}}}{M_{\text{FPR}}} \right) + M_{\text{TPR}} \ln(\epsilon) + M_{\text{FNR}} \ln \left( \frac{M_{\text{FNR}}}{M_{\text{TNR}}} \right) + M_{\text{FNR}} \ln(\epsilon) \right) + M_m \quad (16)$$

$$= -\frac{n}{\ln(2)^2} \left( M_{\text{TPR}} \ln \left( \frac{M_{\text{TPR}}}{M_{\text{FPR}}} \right) + M_{\text{FNR}} \ln \left( \frac{M_{\text{FNR}}}{M_{\text{TNR}}} \right) + \ln(\epsilon) \right) + M_m \quad (17)$$

Dropping constants and constant factors, we get the following cost function

$$C(\tau) = -M_{\text{TPR}} \ln \left( \frac{M_{\text{TPR}}}{M_{\text{FPR}}} \right) - M_{\text{FNR}} \ln \left( \frac{M_{\text{FNR}}}{M_{\text{TNR}}} \right)$$

which we can optimize with dynamic programming in  $\Theta(n \lg n)$  time. We can restrict the choice of  $\tau$  such that

$$\max \left\{ \epsilon \frac{M_{\text{TPR}}}{M_{\text{FPR}}}, \epsilon \frac{M_{\text{FNR}}}{M_{\text{TNR}}} \right\} \leq \epsilon_{\max} \quad (18)$$

and the resulting DBF will have worst-case false positive rate  $\epsilon_{\max}$ .

Here is one sample procedure to optimize the cost function.

1. Sort the learned model confidence scores  $A$  on the training set from least to greatest.
2. For each index  $i$ ,
  - (a) Compute the number of elements with negative class in  $A[:i]$  and with positive class in  $A[i:]$  with a rolling sum.
  - (b) Compute the false/true negative/positive rates of the model given that  $A[:i]$  is classified as negative and  $A[i:]$  is classified as positive.
  - (c) Verify that the optimal backup filter false positive rates are less than  $\epsilon_{\max}$ .
  - (d) Track the index with the lowest cost  $i^*$ .
3. Output the average between  $A[i^* - 1]$  and  $A[i^*]$ .

## 4 Future Work

Rather than outright rejecting decision thresholds  $\tau$ , where the worst-case false positive rate is greater than  $\epsilon_{\max}$ , it might be possible to compute optimal backup filter false positive rates that obey this constraint. It might also be possible to derive optimal parameter values given a fixed memory budget and maximum false positive rate.

## References

- [1] Allison Bishop and Hayder Tirmazi. “Adversary Resilient Learned Bloom Filters”. In: (2025). arXiv: 2409.06556.
- [2] Tim Kraska et al. “The Case for Learned Index Structures”. In: (2018).
- [3] Moni Naor and Yogev Eylon. “Bloom Filters in Adversarial Environments”. In: (2019). DOI: 10.1145/3306193.