

Supplementary Material

1. Improving the First Class of Hashing Algorithms: Maximum Likelihood Estimators

The main goal is to prove **Theorem 1** stated in the main paper. We restate the procedure of maximum likelihood estimation here in this context.

Suppose we have k hashes, and define v_{ik} (equivalently v_{jk}, v_{wk}) to be the respective values under the k^{th} hash. Consider the k^{th} triple given by (v_{ik}, v_{jk}, v_{wk}) . There are only five possible sets of triples: a) all elements are equal, b) two elements are equal and the third is different, and c) all elements are distinct. Table 1 shows the different types of triples.

	A	B	C	D	E
v_{is}	equal	different	equal	equal	distinct
v_{js}	equal	equal	different	equal	distinct
v_{ws}	equal	equal	equal	different	distinct

Table 1: Table of possible triples.

Suppose we count the triples in each set and denote this as n_l , where $l \in \{A, B, C, D, E\}$, and p_l the probability of observing a triple falling in the set l . We note that $n_A + n_B + n_C + n_D + n_E = k$.

Since we have pre-computed $d(\mathbf{x}_i, \mathbf{w})$ and $d(\mathbf{x}_j, \mathbf{w})$, and $f^{(1)}$ is linear, we can invert this function to find the corresponding $\rho_h(\mathbf{x}_i, \mathbf{w})$, $\rho_h(\mathbf{x}_j, \mathbf{w})$. Moreover, we can write $p_A + p_C = \rho_h(\mathbf{x}_i, \mathbf{w})$, $p_A + p_B = \rho_h(\mathbf{x}_j, \mathbf{w})$, $p_A + p_D = \rho_h(\mathbf{x}_i, \mathbf{x}_j)$, and $p_A + p_B + p_C + p_D + p_E = 1$. Finally, since $d(\mathbf{x}_i, \mathbf{x}_j)$ is given by $p_A + p_D$, we write the log likelihood function $\ell(p_A, p_B, p_C, p_D, p_E)$ in terms of p_A and p_D , and get

$$\begin{aligned} \ell(p_A, p_D) = & n_A \log(p_A) + n_B \log(\rho_h(\mathbf{x}_j, \mathbf{w}) - p_A) + n_C \log(\rho_h(\mathbf{x}_i, \mathbf{w}) - p_A) \\ & + n_D \log(p_D) + n_E \log(1 - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w}) + p_A - p_D) \end{aligned} \quad (1)$$

as we want to compute the maximum likelihood estimate of $\hat{p}_A + \hat{p}_D$ to give an estimate of $\rho_h(\mathbf{x}_i, \mathbf{x}_j)$.

Theorem 1 Suppose we have a hashing algorithm where the estimate of interest is given by

$$\rho_h(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}[f(Y^{(1)})] \quad (2)$$

where $Y^{(1)}$ is a Bernoulli random variable, f is a linear function, and the output of the hashing algorithm takes discrete values. Suppose we add a weighted vector \mathbf{w} , and compute

the maximum likelihood estimate via (1). Then: a) this estimator is unbiased, and b) the asymptotic variance of this estimator is always lower than or equal to the variance of the estimate without using the MLE.

Proof In part a), the estimator is given by $\hat{p}_A + \hat{p}_D$. We can see that observe from Table 1 that sets A and D correspond to when the hashed values $v_{is} = v_{js}$ for $1 \leq s \leq k$, and by definition, we have $\rho_h(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{P}[v_i = v_j]$, hence by the law of large numbers, $\sum_{s=1}^k 1_{\{v_{is}=v_{js}\}}/k \approx \rho_h(\mathbf{x}_i, \mathbf{x}_j)$ and hence this estimator is unbiased.

We now prove part b).

Recall that for random variables Y_1, Y_2, \dots, Y_k , we have that

$$\text{Var}\left[\frac{Y_1 + \dots + Y_k}{k}\right] = \frac{1}{k^2} \left(\sum_{s=1}^k \text{Var}[Y_s] + 2 \sum_{s \geq t} \text{Cov}(Y_s, Y_t) \right) \quad (3)$$

In the original case without any modifications, we have $Y_s := 1_{\{v_{is}=v_{js}\}}$ for $1 \leq s \leq k$, and each Y_s are i.i.d. from a Bernoulli distribution.

Hence with k hashes, we must have

$$\text{Var}\left[\frac{Y_1 + \dots + Y_k}{k}\right] = \frac{1}{k^2} \left(\sum_{s=1}^k \text{Var}[Y_s] \right) \quad (4)$$

$$= \frac{\text{Var}[1_{\{v_i=v_j\}}]}{k} \quad (5)$$

$$= \frac{(p_A + p_D)(p_B + p_C + p_E)}{k} \quad (6)$$

Suppose we now look at our new estimator, where we make use of the pre-computed values. From Equation (1), we can compute the partial derivatives to be

$$\frac{\partial \ell}{\partial p_A} = \frac{n_A}{p_A} + \frac{n_B}{p_A - \rho_h(\mathbf{x}_j, \mathbf{w})} + \frac{n_C}{p_A - \rho_h(\mathbf{x}_i, \mathbf{w})} + \frac{n_E}{1 + p_A - p_D - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w})} \quad (7)$$

$$\frac{\partial \ell}{\partial p_D} = \frac{n_D}{p_D} - \frac{n_E}{1 + p_A - p_D - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w})} \quad (8)$$

and Hessian matrix of second partial derivatives as

$$\begin{aligned} H &= \begin{pmatrix} -\frac{n_A}{p_A^2} - \frac{n_B}{(\rho_h(\mathbf{x}_j, \mathbf{w}) - p_A)^2} - \frac{n_C}{(\rho_h(\mathbf{x}_i, \mathbf{w}) - p_A)^2} - \frac{n_E}{(1 - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w}) + p_A - p_D)^2} & \frac{n_E}{(1 - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w}) + p_A - p_D)^2} \\ \frac{n_E}{(1 - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w}) + p_A - p_D)^2} & -\frac{n_E}{(1 - \rho_h(\mathbf{x}_i, \mathbf{w}) - \rho_h(\mathbf{x}_j, \mathbf{w}) + p_A - p_D)^2} - \frac{n_D}{p_D^2} \end{pmatrix} \\ &= \begin{pmatrix} -\frac{n_A}{p_A^2} - \frac{n_B}{p_B^2} - \frac{n_C}{p_C^2} - \frac{n_E}{p_E^2} & -\frac{n_D}{p_D^2} - \frac{n_E}{p_E^2} \\ -\frac{n_D}{p_D^2} - \frac{n_E}{p_E^2} & -\frac{n_D}{p_D^2} - \frac{n_E}{p_E^2} \end{pmatrix} \quad (9) \end{aligned}$$

From the Hessian, we can compute the expected Fisher Information \mathcal{I} of p_A and p_D given by

$$\mathcal{I} = -\mathbb{E}[H] \quad (10)$$

$$= \begin{pmatrix} \frac{\mathbb{E}[n_A]}{p_A^2} + \frac{\mathbb{E}[n_B]}{p_B^2} + \frac{\mathbb{E}[n_C]}{p_C^2} + \frac{\mathbb{E}[n_E]}{p_E^2} & -\frac{\mathbb{E}[n_E]}{p_E^2} \\ -\frac{\mathbb{E}[n_E]}{p_E^2} & \frac{\mathbb{E}[n_D]}{p_D^2} + \frac{\mathbb{E}[n_E]}{p_E^2} \end{pmatrix} \quad (11)$$

$$= \begin{pmatrix} \frac{kp_A}{p_A^2} + \frac{kp_B}{p_B^2} + \frac{kp_C}{p_C^2} + \frac{kp_E}{p_E^2} & -\frac{kp_E}{p_E^2} \\ -\frac{kp_E}{p_E^2} & \frac{kp_D}{p_D^2} + \frac{kp_E}{p_E^2} \end{pmatrix} \quad (12)$$

$$= k \begin{pmatrix} \frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} & -\frac{1}{p_E} \\ -\frac{1}{p_E} & \frac{1}{p_D} + \frac{1}{p_E} \end{pmatrix} \quad (13)$$

Now, $\begin{pmatrix} \hat{p}_A \\ \hat{p}_D \end{pmatrix}$ converges in distribution to a bivariate normal random variable $N\left(\begin{pmatrix} p_A \\ p_D \end{pmatrix}, \mathcal{I}^{-1}\right)$, where

$$\mathcal{I}^{-1} = \frac{1}{k \left[\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) \left(\frac{1}{p_D} + \frac{1}{p_E} \right) - \frac{1}{p_E^2} \right]} \begin{pmatrix} \frac{1}{p_D} + \frac{1}{p_E} & \frac{1}{p_E} \\ \frac{1}{p_E} & \frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \end{pmatrix} \quad (14)$$

We can now write the variance of our new estimator as

$$\text{Var}[\hat{p}_A + \hat{p}_D] = \text{Var}[\hat{p}_A] + \text{Var}[\hat{p}_D] + 2\text{Cov}[\hat{p}_A, \hat{p}_D] \quad (15)$$

$$= \frac{\left(\left(\frac{1}{p_D} + \frac{1}{p_E} \right) + \left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) + \frac{2}{p_E} \right)}{k \left[\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) \left(\frac{1}{p_D} + \frac{1}{p_E} \right) - \frac{1}{p_E^2} \right]} \quad (16)$$

$$= \frac{\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_D} + \frac{4}{p_E}}{k \left[\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) \left(\frac{1}{p_D} + \frac{1}{p_E} \right) - \frac{1}{p_E^2} \right]} \quad (17)$$

Finally, we want to show that the asymptotic variance of this estimator is always lower than or equal to the variance of the original estimator, which is equivalent to showing that

$$(p_A + p_D)(p_B + p_C + p_E) - \frac{\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_D} + \frac{4}{p_E}}{\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) \left(\frac{1}{p_D} + \frac{1}{p_E} \right) - \frac{1}{p_E^2}} \geq 0 \quad (18)$$

Suppose we let

$$A := (p_A + p_D)(p_B + p_C + p_E) \left[\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) \left(\frac{1}{p_D} + \frac{1}{p_E} \right) - \frac{1}{p_E^2} \right] \quad (19)$$

$$B := \frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_D} + \frac{4}{p_E} \quad (20)$$

We then want to show that $A \geq B$ for all values of p_A, p_B, p_C, p_D, p_E .

The key idea here is to replace the "1" terms by $p_A + p_B + p_C + p_D + p_E$ whenever they occur.

Note that A remains unchanged, because we can write

$$A := \frac{(p_A + p_D)}{1} \frac{(p_B + p_C + p_E)}{1} \left[\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E} \right) \left(\frac{1}{p_D} + \frac{1}{p_E} \right) - \frac{1}{p_E^2} \right] \quad (21)$$

and the $p_A + p_B + p_C + p_D + p_E$ terms cancel.

However, B now becomes

$$B := \frac{p_A + p_B + p_C + p_D + p_E}{p_A} + \frac{p_A + p_B + p_C + p_D + p_E}{p_B} + \frac{p_A + p_B + p_C + p_D + p_E}{p_C} \\ + \frac{p_A + p_B + p_C + p_D + p_E}{p_D} + \frac{4(p_A + p_B + p_C + p_D + p_E)}{p_E} \quad (22)$$

Now, we need to show that $A - B \geq 0$ for all $0 \leq p_A, p_B, p_C, p_D, p_E \leq 1$. We first write

$$A - B = \frac{C}{p_A p_B p_C p_D p_E} \quad (23)$$

where

$$C = (p_A - p_D)^2 (p_B + p_C) p_B p_C + (p_B - p_C)^2 (p_A + p_D) p_A p_D \\ + p_A p_E (p_A (p_B + p_C) (p_B + p_C + p_E) - 4 p_B p_C p_D) \quad (24)$$

$$= \alpha p_D^2 + \beta p_D + \gamma \quad (25)$$

as a quadratic in p_D by replacing p_E by $1 - p_A - p_B - p_C - p_D$, with

$$\alpha = (p_A + p_B)(p_A + p_C)(p_B + p_C) \quad (26)$$

$$\beta = 2p_A(p_B p_C(p_B + p_C - 2) + p_A^2(p_B + p_C) + p_A(p_B + p_C - 1)(p_B + p_C)) \quad (27)$$

$$\gamma = p_A^2(p_A + p_B - 1)(p_A + p_C - 1)(p_B + p_C) \quad (28)$$

We can complete the square in our quadratic at (25) and get

$$C = \alpha \left(p_D + \frac{\beta}{2\alpha} \right)^2 - \frac{\beta^2}{4\alpha} + \gamma \quad (29)$$

and now we just need to show that that

$$-\frac{\beta^2}{4\alpha} + \gamma \geq 0 \quad (30)$$

which is a term in p_A, p_B and p_C . The left hand side of (30) gives

$$\frac{p_A^2(p_B - p_C)^2(p_A p_B(1 - p_A - p_B) + p_A p_C(1 - p_A - p_C) + p_B p_C(1 - p_B - p_C) - 2p_A p_B p_C)}{(p_A + p_B)(p_A + p_C)(p_B + p_C)} \quad (31)$$

but we note that since $p_A + p_B + p_C + p_D + p_E = 1$, we can write the third factor of the numerator of (31) as

$$p_A p_B(p_C + p_D + p_E) + p_A p_C(p_B + p_D + p_E) + p_B p_C(p_A + p_D + p_E) - 2p_A p_B p_C \quad (32)$$

and hence note we have a common factor of $p_{APB}p_C$ which cancels out when we expand (32)

$$(p_{APB} + p_{APC} + p_{BPC})(p_D + p_E) + p_{APB}p_C \quad (33)$$

which is always non-negative. We can hence write (31) as

$$\frac{p_A^2(p_B - p_C)^2((p_{APB} + p_{APC} + p_{BPC})(p_D + p_E) + p_{APB}p_C)}{(p_A + p_B)(p_A + p_C)(p_B + p_C)} \geq 0 \quad (34)$$

which implies that (25) is always non-negative, and thus

$$(p_A + p_D)(p_B + p_C + p_E) - \frac{\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_D} + \frac{4}{p_E}}{\left(\frac{1}{p_A} + \frac{1}{p_B} + \frac{1}{p_C} + \frac{1}{p_E}\right)\left(\frac{1}{p_D} + \frac{1}{p_E}\right) - \frac{1}{p_E^2}} \geq 0 \quad (35)$$

with equality when $p_B = p_C$ and $p_D = p_A(1 - p_A - p_B)/(p_A + p_B)$.

This completes the proof. ■

1.1. Random Projections and Very Sparse Random Projections

We prove the following theorem

Theorem 2 Suppose $\mathbf{w}_1, \dots, \mathbf{w}_S$ are chosen to be pairwise orthogonal. With the estimator $v_i v_j$, we can write the control variate

$$Z = \sum_{s,t} \alpha_{s,t} v_{\mathbf{w}_s} v_{\mathbf{w}_t} \quad (36)$$

with $\alpha_{s,t} = d(\mathbf{x}_i, \mathbf{w}_s)d(\mathbf{x}_j, \mathbf{w}_t) + d(\mathbf{x}_i, \mathbf{w}_t)d(\mathbf{x}_j, \mathbf{w}_s)$, with control variate correction $\hat{c} = -\frac{1}{2}$.

Proof For each hash, our estimate of the inner product $\hat{\theta}$ can be written as

$$\hat{\theta} = v_i v_j + \hat{c} \sum_{s,t} \alpha_{s,t} (v_{\mathbf{w}_s} v_{\mathbf{w}_t} - \mathbb{E}[v_{\mathbf{w}_s} v_{\mathbf{w}_t}]) \quad (37)$$

The variance of our estimates can hence be written as

$$\text{Var}[\hat{\theta}] = \text{Var}[v_i v_j + \hat{c} \sum_{s,t} \alpha_{s,t} v_{\mathbf{w}_s} v_{\mathbf{w}_t}] \quad (38)$$

$$\begin{aligned} &= \text{Var}[v_i v_j] + \hat{c}^2 \sum_{s,t} \alpha_{s,t}^2 \text{Var}[v_{\mathbf{w}_s} v_{\mathbf{w}_t}] + 2\hat{c} \sum_{s,t} \alpha_{s,t} \text{Cov}(v_i v_j, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \\ &\quad + 2\hat{c} \sum_{s,t,m,n} \alpha_{s,t} \alpha_{m,n} \text{Cov}(v_{\mathbf{w}_m} v_{\mathbf{w}_n}, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \end{aligned} \quad (39)$$

Suppose we bring in \hat{c} into the α s, i.e. set $\tilde{\alpha} = \hat{c}\alpha$, and we map each of the tuples (s, t) to $k \in \{1, 2, \dots, S(S-1)/2\}$. Then we can rewrite the above equation to get

$$\text{Var}[\hat{\theta}] = \text{Var}[v_i v_j + \sum_k \alpha_k v_{\mathbf{w}_s} v_{\mathbf{w}_t}] \quad (40)$$

$$\begin{aligned} &= \text{Var}[v_i v_j] + \sum_k \alpha_k^2 \text{Var}[v_{\mathbf{w}_s} v_{\mathbf{w}_t}] + 2 \sum_k \alpha_k \text{Cov}(v_i v_j, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \\ &\quad + 2 \sum_{k,k'} \alpha_k \alpha_{k'} \text{Cov}(v_{\mathbf{w}_m} v_{\mathbf{w}_n}, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \end{aligned} \quad (41)$$

Then if we take partial derivatives with respect to α_s , we have for some arbitrary α_k term

$$\begin{aligned} \frac{\partial \text{Var}[\hat{\theta}]}{\partial \alpha_k} &= 2\alpha_k \text{Var}[v_{\mathbf{w}_s} v_{\mathbf{w}_t}] + 2 \text{Cov}(v_i v_j, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \\ &\quad + 2 \sum_{k' \neq k} \alpha_{k'} \text{Cov}(v_{\mathbf{w}_m} v_{\mathbf{w}_n}, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \end{aligned} \quad (42)$$

Since we know the weighted vectors are chosen to be orthogonal to each other, then $\text{Cov}(v_{\mathbf{w}_m} v_{\mathbf{w}_n}, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) = 0$ if $(m, n) \neq (s, t)$ (or $(m, n) \neq (t, s)$). Hence (42) simplifies to

$$\frac{\partial \text{Var}[\hat{\theta}]}{\partial \alpha_k} = 2\alpha_k \text{Var}[v_{\mathbf{w}_s} v_{\mathbf{w}_t}] + 2 \text{Cov}(v_i v_j, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) \quad (43)$$

and equating this to zero, we must have that

$$\alpha_k = -\text{Cov}(v_i, v_j, v_{\mathbf{w}_s} v_{\mathbf{w}_t}) / \text{Var}[v_{\mathbf{w}_s} v_{\mathbf{w}_t}] \quad (44)$$

Since v_i s are distributed multivariate normal, and \mathbf{x}_i s and the weighted vectors are normalized to unit length 1, then any

$$\text{Cov}(v_i v_j, v_k v_l) = d(\mathbf{x}_j, \mathbf{x}_k) d(\mathbf{x}_i, \mathbf{x}_l) + d(\mathbf{x}_j, \mathbf{x}_l) d(\mathbf{x}_i, \mathbf{x}_k) \quad (45)$$

and in general we get

$$\alpha_k = -\frac{1}{2} d(\mathbf{x}_j, \mathbf{x}_k) d(\mathbf{x}_i, \mathbf{x}_l) + d(\mathbf{x}_j, \mathbf{x}_l) d(\mathbf{x}_i, \mathbf{x}_k) \quad (46)$$

Since we brought \hat{c} into the α s, we thus have the optimal control variate coefficient and terms to be

$$\hat{c} = \frac{1}{2} \quad \alpha_{i,j} = d(\mathbf{x}_j, \mathbf{x}_k) d(\mathbf{x}_i, \mathbf{x}_l) + d(\mathbf{x}_j, \mathbf{x}_l) d(\mathbf{x}_i, \mathbf{x}_k) \quad (47)$$

or in general, any form of

$$\hat{c} = \frac{1}{2} \beta \quad \alpha_{i,j} = \frac{1}{\beta} d(\mathbf{x}_j, \mathbf{x}_k) d(\mathbf{x}_i, \mathbf{x}_l) + d(\mathbf{x}_j, \mathbf{x}_l) d(\mathbf{x}_i, \mathbf{x}_k) \quad (48)$$

for any $\beta \in \mathbb{R}$. ■