

Test of diffusion point estimator

Elias Amselem^{§†*},

[†] Department of Cell and Molecular Biology, Science for Life Laboratory, Uppsala University, Husarg. 3, SE-75124 Uppsala, Sweden

Correspondence e-mail: *elias.amselem@icm.uu.se or

[§] These authors contributed equally to this work.

Abstract

Here we describe a set of tests comparing the MSD, CVE and ECVE diffusion estimators for short trajectories.

I. GENERAL COMMENTS

There are many comments and calculations added as lyx Notes in the lyx document, look at them for details on the calculations, a copy of the Lyx file is in the Doc folder in the data folder.

- Data folder: 20220519_EXP-22-BL9418_(SimulationDiffusionTest)
- TriggerFactor_Code: copy of the code, us it for generating data
- Folders: S_wXX_NXXnm (w= window size, N Gaussian noise)
- Folder: S_w16_N70nm_20perodes_2022_09_08 is used for the real-time tracking of Trigger Faktor manuscript.

For theory and results see below.

II. DETERMINING DIFFUSION COEFFICIENTS

Given position estimations for a particle trajectory, we like to find relevant parameters that characterize the underlying dynamics related to it. For dynamical processes where the molecule in question is driven by a diffusion process and associates/dissociates to a larger molecule, it is expected that the diffusion state-change can be observed. A more careful treatment of MSD where contributions from motion blur and position estimation noise are incorporated reveals a richer theory with alternative estimators. Building on the covariance point estimator (CVE) approach, that was derived by Berglund et al. and related work [1–3], we expand their work further to incorporate an arbitrary time lag within the theory and derive alternative point estimators.

Suppose that a particle is moving in one dimension by pure Brownian motion with a diffusion coefficient D and that the particle position X_k for $k = 1, \dots, N$ is estimated with a time interval of Δt , which will be referred to as the frame time, and k is the frame index. For higher dimensions, 2D and 3D motion, we are assuming that each dimension is independent and each can be treated separately. The observed position of the particle is then the average of its position weighted by a “shutter function” $s(t)$, a non-negative function whose integral over the frame time interval is unity. On top of this, we add a noise factor representing the localization noise term which we assume is a zero mean gaussian with covariance $\langle \epsilon_i \epsilon_j \rangle = \sigma^2 \delta_{i,j}$. Given that we have N frames, the k th frame ending at time $t = k \cdot \Delta t$ where $k = 1, \dots, N$, gives an observed position X_k given by

$$X_k = \int_{(k-1)\Delta t}^{k\Delta t} s(t' - (k-1)\Delta t) X_{true}(t') dt' + \epsilon_k \quad (1)$$

where $X_{true}(t')$ is the true position of the particle at time t' and ϵ_k is the value of the additive localization noise in frame k . In the ideal case when $\epsilon_k = 0$ and $s(t)$ is a delta function then the observed position is the true position and the distribution of position intervals $\Delta_{k,j} = X_{k+j} - X_k$ for any fixed $j > 0$ are independent and are zero mean Gaussian distributed with a variance $\langle \Delta_{k,j}^2 \rangle = 2 \cdot D \cdot \Delta t \cdot j$ and the covariance matrix off-diagonal elements $\langle \Delta_{k,j} \Delta_{k',j} \rangle = 0$ for all $k \neq k'$. Introducing localization noise and blur will induce correlations between $\Delta_{k,j}$ which will give a covariance matrix with nonzero off-diagonal elements.

a) *Derivation of the zero mean:* Using eq.(1) it is straightforward to show that $\langle \Delta_{k,j} \rangle = 0$.

b) *Derivation of covariance:* For the covariance we like to evaluate $\langle \Delta_{k,j} \Delta_{k',j} \rangle$, and after some algebra and using the Brownian motion property $\langle X_{true}(t'') X_{true}(t') \rangle = X_{true}(0)^2 + 2D \min(t', t'')$ and that $\min(x, y) = (x + y)/2 - |x - y|/2$ we get the covariance matrix

$$\langle \Delta_{k,j} \Delta_{k',j} \rangle = \begin{cases} 2D(j - 2R)\Delta t + 2\sigma^2, & |k - k'| = 0 \\ 2D(j - |k' - k|)\Delta t, & |k - k'| < j \\ 2DR\Delta t - \sigma^2, & |k - k'| = j \\ 0, & |k - k'| > j \end{cases} \quad (2)$$

where $j > 0$ and

$$R = \frac{1}{\Delta t} \left(\int_0^{\Delta t} s(t')t' \cdot dt' - \int_0^{\Delta t} \int_0^{\Delta t} s(t')s(t'') \cdot \min(t', t'') dt'' dt' \right) = \frac{1}{\Delta t} \int_0^{\Delta t} S(t) (1 - S(t)) dt \quad (3)$$

which is the blur factor with the cumulative shutter function

$$S(t) = \int_0^t s(t') dt'$$

The covariance matrix diagonal elements are now the MSD curve as a function of time lag when viewed as a function of j . The covariance matrix presented in Berglund et al. [1] is recovered when $j = 1$, and it can be seen from eq.(2) that it is a real symmetric band matrix where the band width increases as j increases. Following the CVE method outlined in [4], we solve for D and σ^2 . There are several possibilities for this and we look at three:

$$\hat{D} = \frac{1}{2j\Delta t} \overline{\Delta_{k,j}^2} + \frac{1}{j\Delta t} \overline{\Delta_{k,j} \Delta_{k+j,j}} \quad (4)$$

$$\hat{\sigma}^2 = \frac{R}{j} \overline{\Delta_{k,j}^2} + \left(\frac{2R}{j} - 1 \right) \overline{\Delta_{k,j} \Delta_{k+j,j}}$$

$$\hat{D} = \frac{1}{2(j-n)\Delta t} \overline{\Delta_{k,j} \Delta_{k+n,j}}, \text{ for } n = 1, \dots, j-1 \quad (5)$$

$$\sigma^2 = \frac{1}{2} \overline{\Delta_{k,j}^2} - \frac{(j-2R)}{2(j-n)} \overline{\Delta_{k,j} \Delta_{k+n,j}}$$

$$\hat{D} = \frac{1}{2(j-n)\Delta t} \overline{\Delta_{k,j} \Delta_{k+n,j}}, \text{ for } n = 1, \dots, j-1 \quad (6)$$

$$\hat{\sigma}^2 = \frac{R}{(j-n)} \overline{\Delta_{k,j} \Delta_{k+n,j}} - \overline{\Delta_{k,j} \Delta_{k+j,j}}, \text{ for } n = 1, \dots, j-1$$

here hat refers to estimations and the bar is the mean taken along the diagonal of the covariance matrix, thus $\overline{\Delta_{k,j}^2} = \sum_{k=1}^{N-j} \Delta_{k,j}^2 / (N-j-1) = \langle \Delta_{k,j}^2 \rangle_k$, and $\overline{\Delta_{k,j} \Delta_{k+j,j}} = \langle \Delta_{k,j} \Delta_{k+j,j} \rangle_k$ where we use $\langle \rangle_k$ to indicate the mean over the parameter k . The first pair is constructed by the diagonal matrix elements for $|k - k'| = 0$ and $|k - k'| = j$. This is a direct extension of the CVE [4], and are the same when $j = 1$. For the second and third equations, we note that for $0 < |k - k'| < j$ there is no σ^2 or R dependence and thus D can be estimated directly. For σ^2 there are two choices given this D estimator. In the second pair we use the $|k - k'| = 0$ elements, and in the third pair we use $|k - k'| = j$ elements. The estimator finally used is the second pair with $j = 8$ and $|k - k'| = 1$. For our analysis, we note that for $j > 1$, the off-diagonals between $0 < |k - k'| < j$ are all independent of σ , and a simple estimator for diffusion would be to solve for D and σ^2 for a fixed value of $|k - k'|$. For the first off-diagonal, $|k - k'| = 1$, and a time lag of j , the estimator is

$$\tilde{D}_j = \frac{\langle \Delta_{k,j} \Delta_{k',j} \rangle_1}{2\Delta t \cdot (j-1)} \quad (7)$$

where $\langle \rangle_{|k'-k|}$ is the mean along the $|k' - k| = 1$ off-diagonal of the covariance matrix. We will here refer to this estimator as the extended covariance estimator (ECVE). The localization error, σ , can be estimated by inserting the diffusion estimation \tilde{D}_j into the main diagonal, $|k' - k| = 0$, and solving for σ^2 . The ECVE estimator is compared to the CVE and the MSD in sup.inf.(II-1); the conclusion is that the ECVE is more robust in presence of noise and has less bias compared to the MSD when evaluated over short sections of a trajectory.

c) *Likelihood function* : It is possible to find an inverse to the covariance matrix, eq(2), [5], and by the spectral theorem a real symmetric matrix is diagonalizable by orthogonal matrices. Thus we can argue, as in [3], that there is a orthogonal matrix P such that $P\Sigma P^{-1} = \Lambda$ where $\Sigma_j = \langle \Delta_j \Delta_j^T \rangle$ is the covariance matrix, eq(2), where Δ_j is the column vector with displacement of distance j and Λ is a diagonal matrix with diagonal elements described by the vector λ . Note that this also implies that there is a basis where the measurements Δ_j are decoupled, since the diagonal matrix $\Lambda = P\Sigma_j P^{-1} = P \langle \Delta_j \Delta_j^T \rangle P^T = \langle (P\Delta_j)(P\Delta_j)^T \rangle$. In this basis the measurements $P\Delta_j$ are decoupled and constitute a Brownian motion with step sizes λ . Thus, the probability for each step is given by a Gaussian, and the total log likelihood is the log of the product of all of them. Skipping all constant terms, we get for a given offset j

$$l(\Delta_j) = \ln \left(\prod_i \frac{1}{\sqrt{2\pi\lambda_i}} e^{-\frac{(P\Delta_j)_i^2}{2\lambda_i}} \right) \quad (8)$$

$$= -\frac{1}{2} \ln(2\pi |\Lambda|) - \frac{1}{2} (P\Delta_j)^T \Lambda^{-1} (P\Delta_j)$$

$$= -\frac{1}{2} \ln(2\pi |\Sigma_j|) - \frac{1}{2} \Delta_j^T \Sigma_j^{-1} \Delta_j$$

To include several j offsets in the likelihood, one can consider an interval of j values to use in the likelihood, $1 \leq j \leq J$ which gives

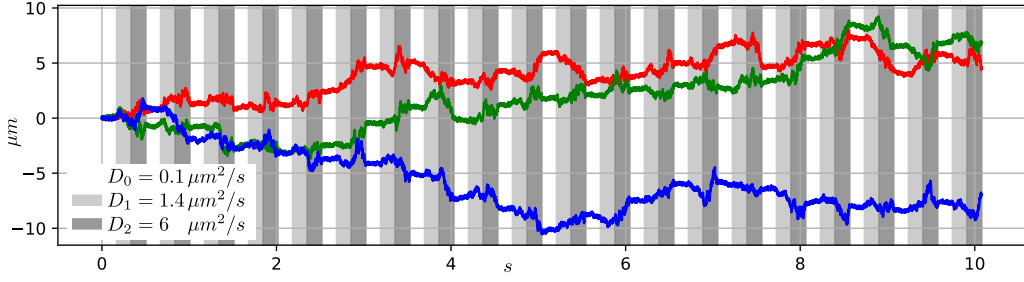


Fig. 1. Trajectory coordinates for x (red), y (green) and z (blue). Each vertical section corresponds to 168ms and within each section a fixed diffusion rate is maintained. The diffusion rates used are $D_0 = 0.1 \mu\text{m}^2/\text{s}$ (white), $D_1 = 1.4 \mu\text{m}^2/\text{s}$ (light gray) and $D_2 = 6 \mu\text{m}^2/\text{s}$ (dark gray)

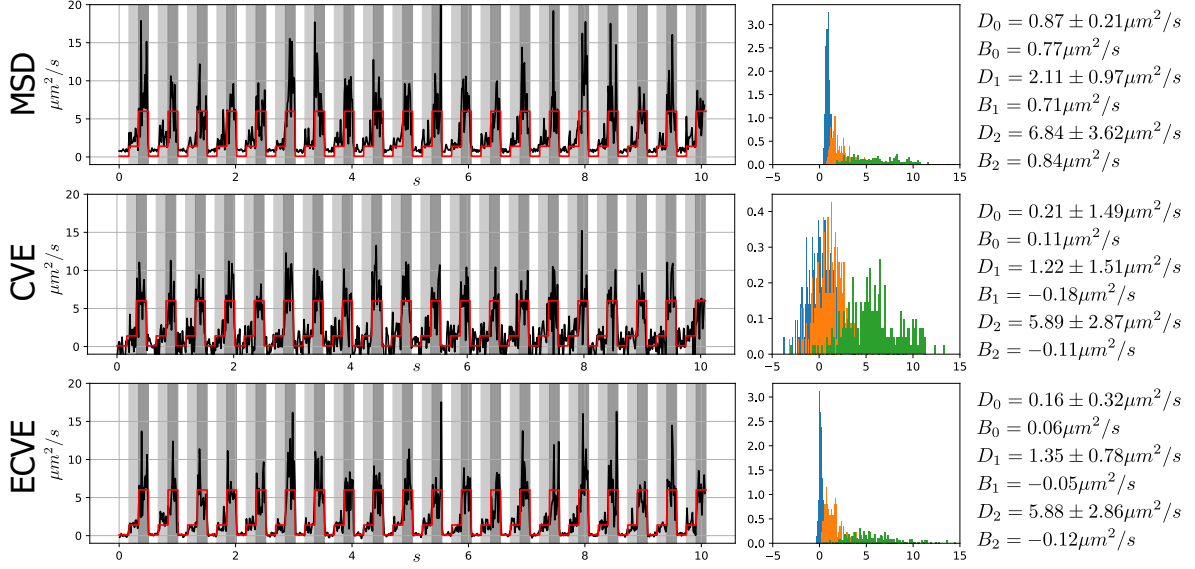


Fig. 2. Estimation of the diffusion over the trajectory in Fig.(1), red is the ground truth. Center plot are the histogram over the trajectory. The color code is by true diffusion; $D_0 = 0.1 \mu\text{m}^2/\text{s}$ (blue), $D_1 = 1.4 \mu\text{m}^2/\text{s}$ (orange) and $D_2 = 6 \mu\text{m}^2/\text{s}$ (green). To the right are for each diffusion rate the mean D_i with error as the distribution standard deviation, and bias B_i between the true value and the mean.

$$L(\{\Delta_j\}_{j=1}^J) = \prod_{j=1}^J \frac{1}{\sqrt{2\pi} |\Sigma_j|} e^{-\frac{1}{2} \Delta_j^T \Sigma_j^{-1} \Delta_j}$$

Since Σ_j depends on D , a maximum likelihood estimator can be obtained by finding D that maximizes the likelihood above. Since this is not a simple estimator, we will not pursue this and instead use the estimator in the main text.

1) *Comparison between MSD, CVE and ECVE diffusion estimations:* Above we discussed the ECVE method and its relation to the mean square displacement (MSD) and the covariance estimator (CVE). To test their performance, we simulate simple 3D diffusion where the diffusion rate is increased in steps ($D_0 = 0.1 \mu\text{m}^2/\text{s}$, $D_1 = 1.4 \mu\text{m}^2/\text{s}$ and $D_2 = 6 \mu\text{m}^2/\text{s}$) where each step is 168ms long. The time step used in the simulation is 0.84ms with a subsampling of 10 samples per time step. On top of the 3D trajectory, we add Gaussian noise with a standard deviation of 70μm per axis. Each axis of the trajectory is shown in Fig.(1) where the colored section corresponds to a diffusion rates. Three point estimators are tested, all are related to the covariance matrix (eq.(2)) and are using a 16 point window for averaging. The first estimator is the MSD which is given by the mean over the diagonal elements of eq.(2) when changing j and applying a linear regression to obtain D and σ^2 . The second estimator is the CVE which have $j = 1$ and we solve for the D , σ^2 and take the mean over the diagonal of the covariance matrix. Lastly is the ECVE which is given by eq.(7) in the main text. In Fig.(2) and Fig.(4) the three tests are shown and compared. The CVE has a poor contrast/large std with the result that the two lower diffusion rates can hardly be distinguished. This is the drawback of using a single step where not enough difference has accumulated. Looking at the MSD, one can identify all three levels but there is a bias for slow diffusion. The ECVE shows better contrast/smaller std compared to both the CVE and MSD and less bias. Estimating the positioning error σ^2 for the three methods (see Fig.(3) and Fig.(4)) shows that the MSD has a large bias that is often negatively valued. Both the CVE and ECVE have better performance and have less bias.

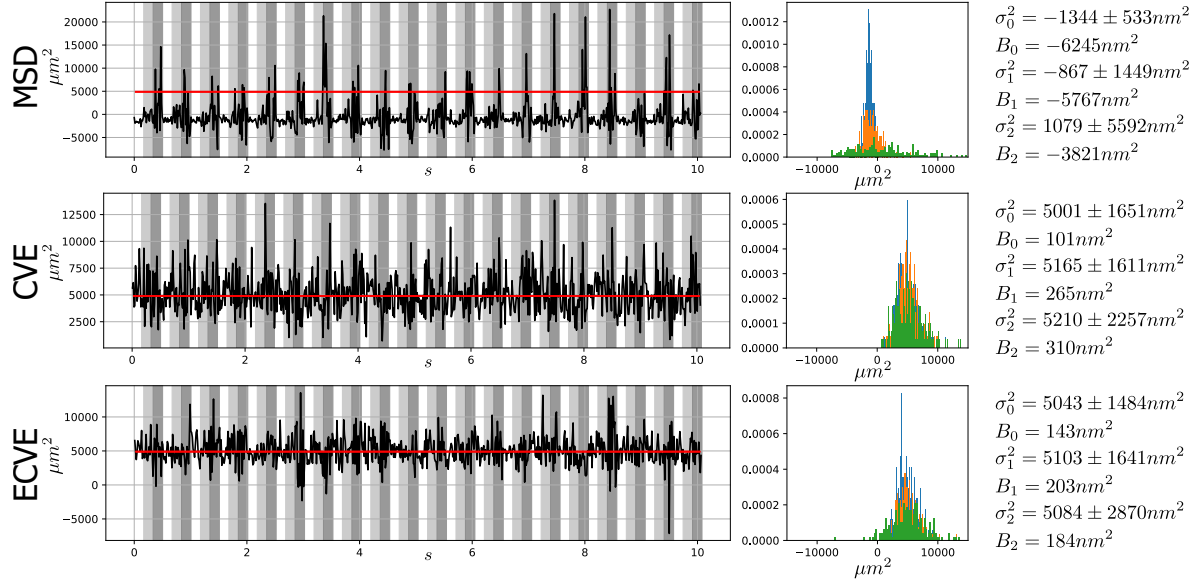


Fig. 3. Estimation of σ^2 the trajectory in Fig.(1), red is the ground truth and corresponds to a standard deviation of $70 \mu\text{m}$. Center, the histogram color coded by true diffusion; $D_0 = 0.1 \mu\text{m}^2/\text{s}$ (blue), $D_1 = 1.4 \mu\text{m}^2/\text{s}$ (orange) and $D_2 = 6 \mu\text{m}^2/\text{s}$ (green). To the right are for each diffusion rate the mean σ_i^2 with error as the distribution standard deviation, and bias B_i between the true value and the mean.

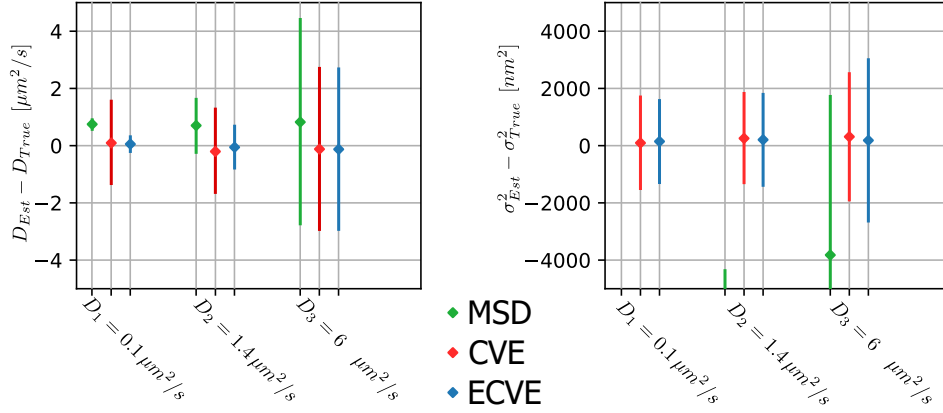


Fig. 4. Comparison between the three diffusion estimators. Left, the difference between the true and estimated diffusion rates is shown for the three diffusion rates tested, the error bars are the std. Right, same as left plot but for the estimated noise σ^2 , here the first MSD point has a large bias and is outside of the plot range.

REFERENCES

1. Berglund, A. J. Statistics of camera-based single-particle tracking. *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics* **82**, 1–8. ISSN: 15393755 (2010).
2. Michalet, X. & Berglund, A. J. Optimal diffusion coefficient estimation in single-particle tracking. *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics* **85**, 1–14. ISSN: 15393755 (2012).
3. Shuang, B. *et al.* Improved analysis for determining diffusion coefficients from short, single-molecule trajectories with photoblinking. *Langmuir* **29**, 228–234. ISSN: 07437463 (2013).
4. Vestergaard, C. L., Blainey, P. C. & Flyvbjerg, H. Optimal estimation of diffusion coefficients from single-particle trajectories. *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics* **89**. ISSN: 15502376 (2014).
5. Allgower, E. L. Exact inverses of certain band matrices. *Numerische Mathematik* **21**, 279–284. ISSN: 09453245 (1973).