

Maximum Entropy Distributions on Unitary Orbits

Polynomial Capacity: Theory, Applications, Generalizations

Jonathan Leake

Technische Universität Berlin

February 4th, 2020

- 1 Last time
 - Maximum entropy distributions
 - Capacity and entropy
- 2 Unitary orbits
 - Orbits and moment polytopes
 - Distributions on orbits and moment polytopes
 - Gelfan-Tsetlin polytope
- 3 Computing max entropy distributions
 - The convex optimization problem
 - The ellipsoid method
- 4 Sampling max entropy distributions
 - Sampling from the associated polytope
 - Returning to the orbit

1 Last time

- Maximum entropy distributions
- Capacity and entropy

2 Unitary orbits

- Orbits and moment polytopes
- Distributions on orbits and moment polytopes
- Gelfan-Tsetlin polytope

3 Computing max entropy distributions

- The convex optimization problem
- The ellipsoid method

4 Sampling max entropy distributions

- Sampling from the associated polytope
- Returning to the orbit

Relative entropy

Given a base measure ν on a support set $\mathcal{S} \subset \mathbb{R}^n$, let $\mu := \phi \cdot \nu$ be another measure given via a function $\phi : \mathcal{S} \rightarrow \mathbb{R}_+$. Define **relative entropy**:

$$D_{\text{KL}}(\mu \parallel \nu) := \int_{\mathcal{S}} \phi(\mathbf{x}) \log \phi(\mathbf{x}) d\nu(\mathbf{x}).$$

Discrete case: Let p be a polynomial associated to μ supported on $\mathcal{S} \subset \mathbb{Z}_+^n$ via $p(\mathbf{x}) = \sum_{\mathbf{v} \in \mathcal{S}} \mu(\mathbf{v}) \cdot \mathbf{x}^{\mathbf{v}}$. When ν is uniform on \mathcal{S} we have

$$D_{\text{KL}}(\mu \parallel \nu) = \sum_{\mathbf{v} \in \mathcal{S}} \left[\frac{\mu(\mathbf{v})}{1/|\mathcal{S}|} \log \frac{\mu(\mathbf{v})}{1/|\mathcal{S}|} \right] \frac{1}{|\mathcal{S}|} = \log |\mathcal{S}| - \mathcal{H}(\mu),$$

where $\mathcal{H}(\mu)$ is the **entropy** of μ . **Properties** for fixed \mathcal{S} :

- When μ is a probability distribution, $\mathcal{H}(\mu) \geq 0$.
- When μ, ν probability distributions, $D_{\text{KL}}(\mu \parallel \nu) \geq 0$.
- $\mathcal{H}(\mu)$ is maximized for $\mu =$ the uniform distribution.
- $D_{\text{KL}}(\mu \parallel \nu)$ is minimized for $\mu = \phi \cdot \nu$ when $\phi \equiv 1$.

Maximum entropy distributions

Last slide: For distributions ν and $\mu := \phi \cdot \nu$ on \mathcal{S} :

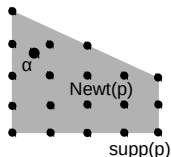
$$D_{\text{KL}}(\mu||\nu) = \int \phi(\mathbf{x}) \log \phi(\mathbf{x}) d\nu(\mathbf{x}).$$

Discrete case: Also $\mathcal{H}(\nu) := - \sum_{\mathbf{v} \in \mathcal{S}} \nu(\mathbf{v}) \log \nu(\mathbf{v})$.

Maximum entropy distribution: Fix α in $\text{hull}(\mathcal{S})$ for discrete \mathcal{S} .

$$\mathcal{H}_{\text{opt}} := \sup_{\substack{\text{supp}(\nu)=\mathcal{S} \\ \mathbb{E}[\nu]=\alpha}} \mathcal{H}(\nu) \quad \text{with} \quad \nu_{\text{opt}} := \arg \sup.$$

Discrete picture for $p \sim \nu$:



Optimum always of the form $\nu_{\text{opt}}(\mathbf{v}) = e^{\langle \mathbf{y}, \mathbf{v} \rangle}$ for some \mathbf{y} (via capacity).

Minimum relative entropy distributions

Previous slide: For distributions ν and $\mu := \phi \cdot \nu$ on \mathcal{S} :

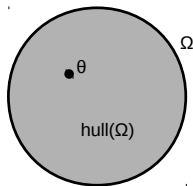
$$D_{\text{KL}}(\mu \parallel \nu) = \int \phi(\mathbf{x}) \log \phi(\mathbf{x}) d\nu(\mathbf{x}).$$

Minimum relative entropy (“maximum entropy”) distribution:

Fix θ in $\text{hull}(\mathcal{S})$ for any \mathcal{S} , and fix some base measure ν on \mathcal{S} .

$$\mathcal{H}_{\text{opt}} := \inf_{\substack{\mu = \phi \cdot \nu \\ \mathbb{E}[\mu] = \theta}} D_{\text{KL}}(\mu \parallel \nu) \quad \text{with} \quad \mu_{\text{opt}} := \arg \inf.$$

Continuous picture where $\Omega = \mathcal{S}$:



Optimum always of the form $\mu_{\text{opt}} = e^{\langle \mathbf{y}, \mathbf{x} \rangle} \cdot \nu$ for some \mathbf{y} (via capacity).

Last slide: Fix θ in $\text{hull}(\mathcal{S})$ for any \mathcal{S} , and base measure ν on \mathcal{S} .

$$\mathcal{H}_{\text{opt}} := \inf_{\substack{\mu = \phi \cdot \nu \\ \mathbb{E}[\mu] = \theta}} D_{\text{KL}}(\mu \| \nu) \quad \text{with} \quad \mu_{\text{opt}} := \arg \inf.$$

This min relative entropy program has the following **dual formulation**:

$$\inf_{\substack{\mu = \phi \cdot \nu \\ \mathbb{E}[\mu] = \theta}} D_{\text{KL}}(\mu \| \nu) = - \inf_{\mathbf{y} \in \mathbb{R}^n} \left[\log \int e^{\langle \mathbf{y}, \mathbf{v} \rangle} d\nu(\mathbf{v}) - \langle \mathbf{y}, \theta \rangle \right] = - \log \text{Cap}_{\theta}(p),$$

where $p(\mathbf{x}) = \int \mathbf{x}^{\mathbf{v}} d\nu(\mathbf{v})$ and $\text{Cap}_{\theta}(p) = \inf_{\mathbf{x} > 0} \frac{\int \mathbf{x}^{\mathbf{v}} d\nu(\mathbf{v})}{\mathbf{x}^{\theta}}$.

Strong duality: $\mu_{\text{opt}} = e^{\langle \mathbf{y}_{\text{opt}}, \mathbf{v} \rangle} \cdot \nu = \mathbf{x}_{\text{opt}}^{\mathbf{v}} \cdot \nu$ where \mathbf{x}_{opt} is the optimum input for $\text{Cap}_{\theta}(p) \implies$ Max entropy optimization equivalent to $\text{Cap}_{\theta}(p)$.

- Infinite-dimensional max entropy vs n -dimensional capacity.
- Ellipsoid method for capacity via evaluation and gradient oracles.

- 1 Last time
 - Maximum entropy distributions
 - Capacity and entropy
- 2 Unitary orbits
 - Orbits and moment polytopes
 - Distributions on orbits and moment polytopes
 - Gelfan-Tsetlin polytope
- 3 Computing max entropy distributions
 - The convex optimization problem
 - The ellipsoid method
- 4 Sampling max entropy distributions
 - Sampling from the associated polytope
 - Returning to the orbit

Unitary group: Set of $n \times n$ complex matrices U such that $U^* = U^{-1}$.

Hermitian matrix: An $n \times n$ complex matrix H such that $H = H^*$.

- $H = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^*$ for $\lambda_1 \geq \dots \geq \lambda_n \in \mathbb{R}$ and \mathbf{v}_i orthonormal.
- **Majorization property:** $\sum_{i=1}^k \lambda_i \geq \sum_{i=1}^k h_{i,i}^\downarrow$ for all k ($=$ for n).
- For any unitary U , matrix UHU^* is Hermitian with same eigenvalues.

Unitary “adjoint” orbits: Orbit of H under conjugation action.

- **Equivalent:** Set of Hermitian matrices with specified eigenvalues.
- **“Adjoint”:** Orbits under action of unitary group on its Lie algebra.

E.g.: Orbit of $D = \text{diag}(1, 0, \dots, 0)$ is set of rank-one projection matrices.

Convex hull of this orbit is the set of trace-1 positive semidefinite matrices.

Majorization: Diagonal entries of any UDU^* gives a simplex element.

Moment map

Last slide: $D = \text{diag}(1, 0, \dots, 0)$, orbit of D :

$$\mathcal{O}_D = \{UDU^* : U \text{ unitary}\} = \{H : \text{eig}(H) = (1, 0, \dots, 0)\}.$$

Fact: The map $\text{diag} : H \mapsto (h_{11}, \dots, h_{nn})$ is such that:

$$\text{diag}(\mathcal{O}_D) = \Delta_n := \left\{ \mathbf{x} \in \mathbb{R}_+^n : \sum_{i=1}^n x_i = 1 \right\} = \text{hull}(\{\sigma \cdot \text{eig}(D) : \sigma \in S_n\}).$$

Proof: For any $\mathbf{x} \in \Delta_n$, pick $H = \sqrt{\mathbf{x}} \cdot \sqrt{\mathbf{x}}^\top$ where $\sqrt{\cdot}$ is entrywise.

More general: For any D , we have that $\text{diag}(\mathcal{O}_D)$ is a convex polytope.

Proof idea: The majorization inequalities translate to:

$$\sum_{i=1}^k \lambda_i \geq \sum_{i=1}^k h_{\sigma(i), \sigma(i)} \quad \text{for all } k, \sigma \quad \text{and} \quad \sum_{i=1}^n \lambda_i = \text{tr}(H) = \sum_{i=1}^n h_{i,i}.$$

These cut out the correct polytope, called the **moment polytope**.

We call the map diag the **moment map**. (From torus action, generally.)

E.g.: $(1, 0, \dots, 0) \rightarrow \text{simplex}$, $(n, n-1, \dots, 1) \rightarrow \text{permutohedron}$.

Measures on orbits

Given real diagonal D , consider unitary conjugation orbit \mathcal{O}_D .

Unitary Haar measure: “Uniform” measure on unitary group $U(n)$.

- Unique unitarily invariant probability measure on $U(n)$:

$$\int_{U(n)} f(U) dU = \int_{U(n)} f(VU) dU = \int_{U(n)} f(UV) dU.$$

- Easily sampled by iteratively constructing random orthonormal basis.
- Pushforward through the map $U \mapsto UDU^*$ gives conjugation-invariant **Haar measure** on the orbit \mathcal{O}_D : $\int_{\mathcal{O}_D} f(H) dH = \int_{\mathcal{O}_D} f(UHU^*) dH$.

Looking ahead: This “uniform” Haar measure is a great candidate for the base measure on \mathcal{O}_D for the maximum entropy program on \mathcal{O}_D .

Further: Push forward this measure through diag: “pushes forward” max entropy distributions. **What’s the base measure on the polytope?**

Last week: For $D = \text{diag}(1, 0, \dots, 0)$, this measure is uniform on Δ_n .

Moment polytope measure for unitary orbits

Last slide: For $D = \text{diag}(1, 0, \dots, 0)$, the pushforward measure through diag is uniform on Δ_n :

$$\int_{\mathcal{O}_D} f(\text{diag}(H)) dH = \int_{\Delta_n} f(\mathbf{x}) d\mathbf{x}.$$

E.g.: Let $Y = \text{diag}(\mathbf{y})$ for some $\mathbf{y} \in \mathbb{R}^n$. Then:

$$\int_{\mathcal{O}_D} e^{\langle Y, H \rangle} dH = \int_{\Delta_n} e^{\langle \mathbf{y}, \mathbf{x} \rangle} d\mathbf{x} \quad \implies \quad \text{max entropy correspondence.}$$

Problem: Does not hold for general D . **Duistermaat-Heckman:**
Pushforward is a piecewise polynomial density times uniform measure.

Solution: Higher-dimensional polytope and another “moment-like” map for which the pushforward is uniform: **Gelfand-Tsetlin (GT) polytope.**

More general question: Can we evaluate the Duistermaat-Heckman measure? Piecewise polynomial density very generally for toric action on symplectic manifold, density is always log-concave, etc...

Gelfand-Tsetlin polytope

Hermitian matrix H with eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$: let \mathcal{R} denote the map $H \mapsto (R_{ij})_{1 \leq i \leq j \leq n}$ where $R_{i,j}$ is the i^{th} largest eigenvalue of the $j \times j$ submatrix in the top left corner of H . **E.g.:** $R_{i,n} = \lambda_i$ and $R_{1,1} = H_{1,1}$.

Cauchy interlacing theorem: $R_{1,j} \geq R_{1,j-1} \geq R_{2,j} \geq R_{2,j-1} \geq R_{3,j} \geq \dots$

Rayleigh triangle: Organizing the $R_{i,j}$ in a natural way:

$$\begin{array}{ccccccccc} R_{1,n} & & R_{2,n} & & R_{3,n} & & R_{4,n} & & \dots & & R_{n,n} \\ & R_{1,n-1} & & R_{2,n-1} & & R_{3,n-1} & & \dots & & R_{n-1,n-1} & \\ & & R_{1,n-2} & & R_{2,n-2} & & \dots & & R_{n-2,n-2} & & \\ & & & & \dots & & & & & & \end{array}$$

Cauchy interlacing theorem \implies the Cauchy inequalities cut out the **Gelfand-Tsetlin polytope** $\text{GT}(\lambda)$ of all triangles with top row λ .

Fun fact: The image $\mathcal{R}(\mathcal{O}_D)$ is $\text{GT}(\lambda)$, and the pushforward of the invariant measure dH on \mathcal{O}_D is the uniform measure on $\text{GT}(\lambda)$.

- 1 Last time
 - Maximum entropy distributions
 - Capacity and entropy
- 2 Unitary orbits
 - Orbits and moment polytopes
 - Distributions on orbits and moment polytopes
 - Gelfan-Tsetlin polytope
- 3 Computing max entropy distributions
 - The convex optimization problem
 - The ellipsoid method
- 4 Sampling max entropy distributions
 - Sampling from the associated polytope
 - Returning to the orbit

Ellipsoid method for convex optimization: Iterations of the algorithm:

- 1 Have an ellipsoid which contains the optimum input.
- 2 Compute the gradient at the center of the ellipsoid: Determines which half of the ellipsoid contains the optimum.
- 3 Construct a new ellipsoid containing that half.

At the end: Small ellipsoid containing the optimum \implies approximation.

Need: (1) Oracle for the gradient, and (2) starting ellipsoid.

(There are more details to this, but this is roughly the idea. We'll discuss this a bit further, but see [L-Vishnoi '20], [Singh-Vishnoi '15] for details.)

Recall the convex optimization (capacity) problem:

$$-\inf_{\mathbf{y} \in \mathbb{R}^n} \left[\log \int e^{\langle \mathbf{y}, \mathbf{v} \rangle} d\nu(\mathbf{v}) - \langle \mathbf{y}, \boldsymbol{\theta} \rangle \right] = -\inf_{\mathbf{y} \in \mathbb{R}^n} \left[\log \int e^{\langle \mathbf{y}, \mathbf{v} - \boldsymbol{\theta} \rangle} d\nu(\mathbf{v}) \right].$$

Need: Oracle for $\nabla \log \int e^{\langle \mathbf{y}, \mathbf{v} \rangle} d\nu(\mathbf{v})$ and “bounding box” for \mathbf{y}_{opt} .

Applying ellipsoid method to unitary orbit capacity

Unitary orbits: \mathcal{O}_D for real diagonal D with invariant Haar measure dH :

$$\text{Optimize} \quad - \inf_Y \left[\log \int_{\mathcal{O}_D} e^{\langle Y, H - \Theta \rangle} dH \right] \quad \text{for } \Theta \in \text{hull}(\mathcal{O}_D).$$

Last slide: Need oracle for $\nabla \log \int e^{\langle Y, H \rangle} dH$ and bound on Y_{opt} .

First: dH is a pushforward measure of dU through $U \mapsto UDU^*$:

$$\int_{\mathcal{O}_D} e^{\langle Y, H \rangle} dH = \int_{U(n)} e^{\langle Y, UDU^* \rangle} dU.$$

HCIZ formula: For $\text{eig}(D) = (\lambda_1, \dots, \lambda_n)$ and $\text{eig}(Y) = (y_1, \dots, y_n)$:

$$\int_{U(n)} e^{\langle Y, UDU^* \rangle} dU = \frac{\prod_{i=1}^{n-1} i!}{\prod_{i < j} (\lambda_i - \lambda_j) \prod_{i < j} (y_i - y_j)} \cdot \det \left(\left[e^{\lambda_i y_j} \right]_{i,j=1}^n \right).$$

Apply gradient to get the oracle. **Also:** L'Hopital to handle multiplicity.

Next: Bound on Y_{opt} . Unitary invariance \implies “balanced” measure.

Depends on distance δ of Θ from boundary: $\|Y_{\text{opt}}\| \leq \text{poly}(n, \log \|D\|, \frac{1}{\delta})$.

Alternative method for computing max entropy

Previous slides: Invariant measure dH on \mathcal{O}_D corresponds to uniform measure on a polytope (either through diag of Raleyigh map \mathcal{R}).

Rank-one case: Polytope is simplex and $Y = \text{diag}(\mathbf{y})$ gives

$$\int_{\mathcal{O}_D} e^{\langle Y, H \rangle} dH = \int_{\Delta_n} e^{\langle \mathbf{y}, \mathbf{x} \rangle} d\mathbf{x}.$$

Therefore: Computing integrals / max entropy over \mathcal{O}_D is equivalent to integrals / max entropy over Δ_n . Utilize algos which compute integrals over convex polytopes (e.g. [Lovász-Vempala '06]).

Similarly: For general D and $Y = \text{diag}(\mathbf{y})$, we have

$$\int_{\mathcal{O}_D} e^{\langle Y, H \rangle} dH = \int_{\text{GT}(\lambda)} e^{\langle \mathbf{y}, \text{type}(R) \rangle} dR.$$

Here: R is a Rayleigh triangle, and $\text{type}(R)$ outputs the vector of “diagonal entries” of R (i.e., $h_{kk} = \sum_{i=1}^k R_{i,k} - \sum_{i=1}^{k-1} R_{i,k-1}$).

Again: Utilize algos which compute integrals over convex polytopes.

- 1 Last time
 - Maximum entropy distributions
 - Capacity and entropy
- 2 Unitary orbits
 - Orbits and moment polytopes
 - Distributions on orbits and moment polytopes
 - Gelfan-Tsetlin polytope
- 3 Computing max entropy distributions
 - The convex optimization problem
 - The ellipsoid method
- 4 Sampling max entropy distributions
 - Sampling from the associated polytope
 - Returning to the orbit

Sampling from the associated polytope

Recall: The uniform distribution (dH) on \mathcal{O}_D is easy to sample from.

Idea: Choose random unitary and conjugate D to get sample UDU^* .

Problem: What about max entropy distributions $e^{\langle Y, H \rangle} dH$?

Idea: Sample from the polytope, then “bring back” to the orbit.

Previous slides: For $D_0 = (1, 0, \dots, 0)$, general D , $Y = \text{diag}(\mathbf{y})$:

$$\int_{\mathcal{O}_{D_0}} e^{\langle Y, H \rangle} dH = \int_{\Delta_n} e^{\langle \mathbf{y}, \mathbf{x} \rangle} d\mathbf{x}, \quad \int_{\mathcal{O}_D} e^{\langle Y, H \rangle} dH = \int_{\text{GT}(\lambda)} e^{\langle \mathbf{y}, \text{type}(R) \rangle} dR.$$

Ways to sample: General methods for sampling from log-concave (even log-linear) distributions from convex polytopes.

Or: One can get a good handle on the entrywise CDFs conditioned on various other coordinates: always determinants here. **Sample?**

Next problem: How to go back to the orbit \mathcal{O}_D ?

Returning to the orbit: Rank-one case

Last slide: Sample \mathbf{x} from Δ_n according to $e^{\langle \mathbf{y}, \mathbf{x} \rangle} d\mathbf{x}$.

Next question: How do we convert \mathbf{x} into a sample from \mathcal{O}_D (the orbit of rank-one projections) according to $e^{\langle Y, H \rangle} dH = e^{\langle Y, \mathbf{v}\mathbf{v}^* \rangle} d\mathbf{v}$?

Let's compute the fiber of a given $\mathbf{x} \in \Delta_n$:

$$\text{diag}^{-1}(\mathbf{x}) = \{ \mathbf{v}\mathbf{v}^* : \text{diag}(\mathbf{v}\mathbf{v}^*) = \mathbf{x} \} = \mathbb{T}^n \cdot \sqrt{\mathbf{x}}\sqrt{\mathbf{x}}^\top.$$

That is: $\mathbf{v} = (e^{i\theta_1}\sqrt{x_1}, \dots, e^{i\theta_n}\sqrt{x_n})$. Since Y is diagonal, we then have

$$e^{\langle Y, \mathbf{v}\mathbf{v}^* \rangle} = e^{\langle \text{diag}(e^{-i\theta})Y \text{diag}(e^{i\theta}), \sqrt{\mathbf{x}}\sqrt{\mathbf{x}}^\top \rangle} = e^{\langle Y, \sqrt{\mathbf{x}}\sqrt{\mathbf{x}}^\top \rangle}$$

That is: The max entropy distribution is uniform on fibers.

Therefore: To construct a rank-1 projection sample from a simplex sample, we just need to uniformly sample from \mathbb{T}^n and multiply by $\sqrt{\mathbf{x}}$.

Returning to the orbit: General case

Now: Sample Rayleigh triangle R from $GT(\lambda)$ according to $e^{\langle y, \text{type}(R) \rangle} dR$.

Next: How to convert R into a sample from \mathcal{O}_D according to $e^{\langle Y, H \rangle} dH$?

Problem: Fibers of R are more complicated. (Recall that R holds the eigenvalues of all leading principal submatrices).

Solution: Induct on the rows of R . Assume we have $k \times k$ sample according to the bottom k rows, and construct a $(k+1) \times (k+1)$ matrix using the next $((k+1)^{\text{st}})$ row up.

$$\begin{array}{cccccc} R_{1,n} & R_{2,n} & R_{3,n} & R_{4,n} & \cdots & R_{n,n} \\ & R_{1,n-1} & R_{2,n-1} & R_{3,n-1} & \cdots & R_{n-1,n-1} \\ & & R_{1,n-2} & R_{2,n-2} & \cdots & R_{n-2,n-2} \\ & & & & \ddots & \end{array}$$

That is: Given $k \times k$ matrix X_0 , we want to sample $\begin{bmatrix} X_0 & u \\ u^* & c \end{bmatrix}$.

Key: Max entropy distributions are uniform on these “partial fibers”.

Returning to the orbit: Partial fibers

Task: Given $k \times k$ Hermitian X_0 with eigenvalues $\lambda_1 \geq \dots \geq \lambda_k$, want to sample uniformly random matrix $\begin{bmatrix} X_0 & \mathbf{u} \\ \mathbf{u}^* & c \end{bmatrix}$ with eigenvalues $\gamma_1 \geq \dots \geq \gamma_{k+1}$.

For simplicity: We assume entries of γ and λ are all distinct.

- 1 Diagonalize X_0 using $\begin{bmatrix} U & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{bmatrix}$ to get $\begin{bmatrix} D_0 & \mathbf{w} \\ \mathbf{w}^* & c \end{bmatrix}$ where $\text{diag}(D_0) = \lambda$.
- 2 Compute $c = \sum_{i=1}^{k+1} \gamma_i - \sum_{i=1}^k \lambda_i$.
- 3 **Hardest part:** Sample \mathbf{w} . Recall that λ, γ interlace and note that:

$$\begin{aligned} \prod_{i=1}^{k+1} (z - \gamma_i) &= \det \left(zI - \begin{bmatrix} D_0 & \mathbf{w} \\ \mathbf{w}^* & c \end{bmatrix} \right) \\ &= (z - c) \prod_{i=1}^k (z - \lambda_i) - \sum_{i=1}^k |w_i|^2 \prod_{j \neq i} (z - \lambda_j). \end{aligned}$$

Plug in $z = \lambda_j$, rearrange: $|w_j|^2 = \frac{\prod_{i=1}^{k+1} (\lambda_j - \gamma_i)}{\prod_{i \neq j} (\lambda_j - \lambda_i)}$, sample $e^{i\theta} |w_j|$.

- 4 Apply $\begin{bmatrix} U & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{bmatrix}^*$ to convert D_0 back to X_0 .