

Review of Linear Algebra.

1. Multiplication.

x is a column vector. $x^T y$ is scalar xy^T is a matrix
 x^T is a row vector

$$AB = A[b_1 \dots b_n] = [Ab_1 \dots Ab_n]$$

$$AB = \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} B = \begin{bmatrix} a_1^T B \\ \vdots \\ a_n^T B \end{bmatrix}$$

$$AB = (a_1 \dots a_n) \begin{bmatrix} b_1^T \\ \vdots \\ b_n^T \end{bmatrix} = \sum_{i=1}^n a_i b_i^T \quad \text{sum of matrix}$$

$$AB = \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} [b_1 \dots b_n] = \begin{bmatrix} a_1^T b_1 & a_1^T b_2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \quad \text{matrix of scalar}$$

$$U \Lambda U^T \text{ (when } \Lambda \text{ is diagonal)} = [u_1 \dots u_n] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix} = \begin{bmatrix} \lambda_1 u_1 u_1^T & & \\ & \ddots & \\ & & \lambda_n u_n u_n^T \end{bmatrix}$$

$$= [\lambda_1 u_1 \dots \lambda_n u_n] \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix} = \sum_i \lambda_i u_i u_i^T$$

2. Concepts to describe matrix.

① Hermitian / Symmetric: $A^H = A$, $A^T = A$

② Unitary / Orthogonal: $U^H U = U U^H = I$, $U^T U = U U^T = I$

Orthogonal matrix have some nice properties:

$$① U^T U = \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix} [u_1 \dots u_n] = \begin{bmatrix} u_1^T u_1 & u_1^T u_2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

cols of U form a basis

$$② \text{ norm preserving: } \|u x - u y\|^2 = (u x - u y)^T (u x - u y) = (x - y)^T u^T u (x - y) = \|x - y\|^2$$

③ PSD / PD.

matrix A is PD if $\forall x \neq 0, x^T A x > 0$.

PSD if $\forall x \neq 0, x^T A x \geq 0$

3. Symmetric matrix have nice properties. They have special decomposition and their eigenvalues are indicative of other properties.

Spectral Thm: If A is Hermitian, then $A = U \Lambda U^H$ where U is unitary and Λ is real diagonal.

If A is real and $A = A^T$ (Symmetric), then $A = U \Lambda U^T$ where U is orthonormal.

This is actually eigen-decomp of A :

$$A = U \Lambda U^T \Rightarrow A U = U \Lambda \Rightarrow A [u_1 \dots u_n] = [\lambda_1 u_1 \dots \lambda_n u_n]$$

$$\Rightarrow A u_i = \lambda_i u_i. \quad \lambda_i \text{ are eigenvalues.}$$

② For symmetric matrix, its eigenvalues are indicative of PD.

$$A \text{ is Hermitian/Symmetric} \Rightarrow \begin{cases} A \text{ is PD} \Leftrightarrow \lambda_i > 0 \\ A \text{ is PSD} \Leftrightarrow \lambda_i \geq 0 \end{cases}$$

$$\text{note } x^T A x = \sum_{i=1}^n \lambda_i x^T u_i u_i^T x = \sum_{i=1}^n \lambda_i \|u_i^T x\|^2$$

4. PD is strongly related to symmetry. Sometimes, we just define PD in terms of symmetric matrix

① A matrix can be PD but not symmetric. $\begin{bmatrix} 1 & i \\ i & 1 \end{bmatrix}$

but its eigenvalues are now complex: $\lambda_i = 1 \pm i$
 we don't like this!

② most time, we care PD in terms of symmetric matrix

5. example.


For Gaussian, Σ is symmetric and PD.

① symmetric $\Rightarrow \Sigma = U \Lambda U^T$.


② PD $\Rightarrow \lambda_i > 0$.

so $(x-\mu)^T \Sigma^{-1} (x-\mu) = (x-\mu)^T U \Lambda^{-1} U^T (x-\mu) = (x'-\mu')^T \Lambda^{-1} (x'-\mu')$
 $= \sum \frac{(x'_i - \mu'_i)^2}{\lambda_i}$ since $\lambda_i > 0$, result is an ellipse!

6. Orthogonal and Independence. (note we talk about vectors here)

 v_1 and v_3 are orthogonal. also independent.
 v_2 and v_3 are independent but not orthogonal.
 \Rightarrow independence is a weaker condition.
 v_2 and v_3 are independent, iff v_2 has an orthogonal component w.r.t. v_3 . note v_2 could be decomposed, and one component is $\perp v_3$.

7. Intuition of Subspace.

- ① A subspace is a box. it ~~contains~~ ^{can be spanned} by some independent vectors 
 - ② A matrix can be seen as a collection of vectors $[a_1 \dots a_p]$ thus has a correspondence to a subspace. whether its col vectors are linearly independent is an interesting property.
 - ③ If $a_1 \dots a_p$ are linearly independent. let $\langle A \rangle$ denote subspace spanned. $\dim(\langle A \rangle) = \text{rank}(A) = p$. denote $a_i \in \mathbb{R}^N$. thus $p \leq N$.
 $\dim(\langle A \rangle^\perp) = N - p$.
- Intuition: A N -length vector lives in N -dim space, or can be projected to space whose $\dim < N$. But cannot be lifted to space whose $\dim > N$. so if $a_1 \dots a_p$ are independent, then $p \leq N$.
- ④ Space $\langle A \rangle$ and $\langle A \rangle^\perp$ form orthogonal decomposition of space \mathbb{R}^N . Any vector in $\langle A \rangle$ is orthogonal to every vector in $\langle A \rangle^\perp$.

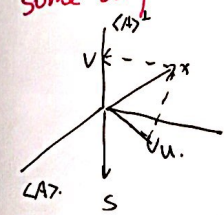
8. Signal as a vector

① A signal can be seen as a vector
 e.g. $x_n = 0.2n^2 + 0.1n + 0.0 = [0.1, n, n^2] \begin{bmatrix} 0.0 \\ 0.1 \\ 0.2 \end{bmatrix} \in \mathbb{R}^3$.

② A bunch of signals live in a subspace:
 $s(1) \dots s(n)$ lives in \mathbb{R}^3

③ Assumption.
 we assume $x = s + w$. s is the signal. w is noise and x is what we observed.

we assume s lives in space $\langle A \rangle$. w must have some component in $\langle A \rangle^\perp$.



$x = u + v = s + w$
 w may not totally live in $\langle A \rangle^\perp$

idea: if we can project x into $\langle A \rangle$. we can have u , thus some sense of s . But how to project a vector on to a subspace?

q. Projection is denoted by an operator π_A .

$\pi_A x = u$. where $\langle A \rangle$ is the subspace to project on.

① matrix A can express π_A : $\pi_A = A(A^H A)^{-1} A^H$.
 note: let $u \in \langle A \rangle$, $v \in \langle A \rangle^\perp = \langle B \rangle$. then $u = A\theta$, $v = B\phi$.
 $\pi_A x = A(A^H A)^{-1} A^H A\theta + A(A^H A)^{-1} A^H B\phi$
 $= A\theta = u$.

Review of MVG.

1. Covariance.

① Covariance is a measure of joint variability of two random variables.

e.g. if X has large value, then Y is likely to have large value. then X and Y have positive covariance.

② $X = [x_1 \dots x_n]^T$ where each component x_i is a random variable, we call X a random vector.

③ Random vector has covariance matrix: pairwise covariance of components.

$$\begin{bmatrix} \text{cov}(x_1, x_1) & \dots & \text{cov}(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \text{cov}(x_n, x_1) & \dots & \text{cov}(x_n, x_n) \end{bmatrix} \quad \text{calculation} \quad \text{cov}(X, X) = E[(X - \mu)(X - \mu)^T] = \Sigma$$

④ $\text{cov}(X, X) = \text{Var}(X)$. For one r. variable, its covariance with itself is called Variance, which is a characteristic of uncertainty.

2. MVG. Definition.

$$\phi(x) = \phi(x; \mu, \Sigma) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$$

where Σ is symmetric and p.d.

① MVG is fully characterized by μ, Σ . denote $X \sim N(\mu, \Sigma)$

② Properties of MVG can be proven by its characteristic function.

characteristic function:

$$\Phi_X(\omega) = E[e^{j\omega^T X}] = \int e^{j\omega^T x} f(x) dx = e^{-j\omega^T \mu - \frac{1}{2}\omega^T \Sigma \omega}$$

3. Properties of MVG.

① Linear transformation.

$$X \sim N(\mu, \Sigma)$$

$$Y = AX$$

$$\Rightarrow Y \sim N(A\mu, A\Sigma A^T)$$

$$\text{note } E_Y = E[e^{j\omega^T Y}] = E[e^{j\omega^T AX}] = e^{j\omega^T A\mu - \frac{1}{2}\omega^T A\Sigma A^T \omega}$$

linear transform of Gaussian is Gaussian. $\sim N(A\mu, A\Sigma A^T)$

e.g. $\frac{1}{n} \sum_{i=1}^n x_i$ is Gaussian, where x_i is Gaussian

$$\text{e.g. } Y = AX + b \Rightarrow Y \sim N(A\mu + b, A\Sigma A^T)$$

many things can be expressed as linear transformation

② Marginals.

$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$\text{then } x_1 \sim N(\mu_1, \Sigma_{11}) \quad x_2 \sim N(\mu_2, \Sigma_{22})$$

$$\text{note } x_1 = [1, 0] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = AX$$

③ Conditioning.

$$x_2 | x_1 = x_1 \sim N(\tilde{\mu}, \tilde{\Sigma}) \quad \tilde{\mu} = \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (x_1 - \mu_1)$$

$$\tilde{\Sigma} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$$

To prove this, we present a different question.

$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}. \quad x_1 | x_2 = x_2 \sim N(E(x_1 | x_2 = x_2), C_{x_1 | x_2}).$$

$$\text{where } C_{x_1 | x_2} = C_{x_1} - A C_{x_2} A^T, \quad E(x_1 | x_2 = x_2) = A(x_2 - \mu_2) + \mu_1.$$

$$\text{and } A \text{ solves } A C_{x_2} = C_{x_1 x_2}$$

$$(a). (x_1 - \mu_1) - A(x_2 - \mu_2) \text{ and } x_2 \text{ are uncorrelated.}$$

$$E[(x_1 - \mu_1) - A(x_2 - \mu_2) \cdot x_2^T] = C_{x_1 x_2} - A C_{x_2} = 0$$

(b1. ~~$x_1 = Ax_2$~~)

$$\begin{bmatrix} (x_1 - \mu_1) - A(x_2 - \mu_2) \\ x_2 \end{bmatrix} = \begin{bmatrix} I & -A \\ 0 & I \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} A\mu_2 - \mu_1 \\ 0 \end{bmatrix}$$

\Rightarrow is Gaussian...

so $(x_1 - \mu_1) - A(x_2 - \mu_2)$ and x_2 are independent.

(c) $\Phi[e^{u'x_1 | x_2 = x_2}] = e^{\underbrace{u'(\mu_1 + A(x_2 - \mu_2))}_{\mu'} - \underbrace{u' C C^T A^T A x_2}_{\Sigma'} / 2}$

solving μ', Σ' gives you the result.

Given x is Gaussian \Rightarrow Marginal are Gaussian.
conditional

what about the other direction?

4. Affine Transformation.

Given $p(x_a) = \mathcal{N}(x_a; \mu_a, \Sigma_a)$

$p(x_b | x_a) = \mathcal{N}(x_b; Mx_a + b, \Sigma_{b|a})$

$\Rightarrow p(x_a, x_b) = \mathcal{N}\left(\begin{bmatrix} x_a \\ x_b \end{bmatrix}; \begin{bmatrix} \mu_a \\ M\mu_a + b \end{bmatrix}, R\right)$

$$R = \begin{pmatrix} \Sigma_a & \Sigma_a M^T \\ M \Sigma_a & \Sigma_{b|a} + M \Sigma_a M^T \end{pmatrix}$$

$\Rightarrow p(x_b) = \mathcal{N}(x_b; \mu_b, \Sigma_b)$ where $\mu_b = M\mu_a + b$
 $\Sigma_b = \Sigma_{b|a} + M \Sigma_a M^T$

$$\Rightarrow p(x_a | x_b) = \mathcal{N}(x_a; \mu_{a|b}, \Sigma_{a|b})$$

very complex but calculable

Summary: $p(x_a, x_b) \Rightarrow p(x_a) p(x_b | x_a)$

$$\frac{p(x_a)}{p(x_b | x_a)} \nRightarrow p(x_a, x_b) \Rightarrow p(x_b), p(x_a | x_b)$$

Everything is a Gaussian! Affine preserve Gaussian

4. Conditional Probability.

conditional probability is also a probability. We can view it as a new probability.

e.g. $p(x_1 | x_2, x_3)$ if we let $\tilde{p}(x)$ to denote $p(x | x_3)$.

$$\text{then } p(x_1 | x_2, x_3) = \tilde{p}(x_1 | x_2) = \frac{\tilde{p}(x_1 | x_2) \tilde{p}(x_2)}{\tilde{p}(x_2)} = \frac{p(x_1 | x_2, x_3) p(x_2 | x_3)}{p(x_2 | x_3)}$$

5. A Linear System Example.

① Given $x_{t+1} = (Ax_t + b_t) + w_t$ where b_t, d_t are known vector.
 $y_t = (Cx_t + d_t) + e_t$ $w_t \sim \mathcal{N}(0, Q)$ $e_t \sim \mathcal{N}(0, R)$

we also know $x_1 \sim \mathcal{N}(\bar{x}_1 | p_1, P_{10})$.

Here we use a notation: $p(x_t | y_{1:t-1}) = \mathcal{N}(x_t; \hat{x}_{t|t-1}, P_{t|t-1})$
when $t=1$ $\mathcal{N}(\bar{x}_1 | p_1, P_{10})$.

② we know x_1 , we can observe y_t , we want to know x_t .

so the problem is $p(x_t | y_{1:t}) \rightarrow p(x_{t+1} | y_{1:t})$.

③ First of all everything is linear transform of Gaussian. Everything is Gaussian.

④ denote $p(x_t | y_{1:t-1}) = \mathcal{N}(\hat{x}_{t|t-1}, P_{t|t-1})$ $\tilde{p}(x_t) = p(\cdot | y_{1:t-1})$

⑤ $p(y_t | x_t, y_{1:t-1})$. since $y_t = (Cx_t + dt) + e_t$, when x_t is given

$y_t | x_t \sim \mathcal{N}(Cx_t + dt, R)$. note $Cx_t + dt$ is given e_t is random variable: $e_t + \text{const.}$
since x_t is observed now, so observing $y_{1:t-1}$ does not matter.

$$\Rightarrow p(y_t | x_t, y_{1:t-1}) = \mathcal{N}(Cx_t + dt, R) = \tilde{p}(y_t | x_t)$$

⑥ $\tilde{p}(x_t) \not\Rightarrow \tilde{p}(x_t | y_t)$ is known. $\Rightarrow p(x_t | y_{1:t})$ is known.
 $\tilde{p}(y_t | x_t)$

denote $\tilde{p}(x_t | y_t) = \mathcal{N}(\hat{x}_{t|t}, P_{t|t})$.

then $\hat{x}_{t|t}, P_{t|t}$ can be derived by $\hat{x}_{t|t-1}, P_{t|t-1}, C, dt, R$

⑦ denote $\tilde{p}(x_t | y_t) = \tilde{\tilde{p}}(x_t)$. $\tilde{\tilde{p}}(x_t) = \tilde{p}(\cdot | y_t) = p(\cdot | y_{1:t})$

⑧ $p(x_{t+1} | x_t, y_{1:t})$ since x_t is directly observed, $y_{1:t}$ does not matter.

$$= p(x_{t+1} | x_t) = \mathcal{N}(Ax_t + b_t, Q) = \tilde{\tilde{p}}(x_{t+1} | x_t)$$

⑨ $\tilde{\tilde{p}}(x_t) \xrightarrow{\text{known}} \tilde{\tilde{p}}(x_{t+1})$ is known $\Rightarrow p(x_{t+1} | y_{1:t})$ is known
 $\tilde{\tilde{p}}(x_{t+1} | x_t)$

Recall our Affine transform:

$$\begin{aligned} \hat{x}_{t+1} &= A\hat{x}_{t|t} + b_t \\ P_{t+1|t} &= AP_{t|t}A^T + Q. \end{aligned}$$

By knowing x_t , observe $y_{1:t}$. we can say something about x_{t+1} .

6. Canonical form

① Gaussian is fully characterized by $\mu, \Sigma \rightarrow$ its first and second moment. This is moment form

② Canonical form: introduce information vector ξ and information matrix Ω .

$$\begin{aligned} \xi &\leftrightarrow \mu \\ \Omega &\leftrightarrow \Sigma \end{aligned} \quad \begin{cases} \Omega = \Sigma^{-1} \\ \Sigma = \Omega^{-1} \end{cases} \quad \begin{cases} \xi = \Sigma^{-1}\mu \\ \mu = \Omega^{-1}\xi \end{cases}$$

③ derive:

$$\begin{aligned} p(x) &= \det(2\pi\Sigma)^{-\frac{1}{2}} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)) \\ &= \eta \exp(-\frac{1}{2}x^T \Sigma^{-1}x + x^T \Sigma^{-1}\mu - \frac{1}{2}\mu^T \Sigma^{-1}\mu) \\ &= \eta' \exp(-\frac{1}{2}x^T \Omega x + x^T \xi) \end{aligned}$$

$$\xi = \begin{bmatrix} \xi_0 \\ \xi_1 \end{bmatrix} \quad \Lambda = \begin{bmatrix} \Lambda_{00} & \Lambda_{01} \\ \Lambda_{10} & \Lambda_{11} \end{bmatrix}$$

④ Marginalization

$$\begin{aligned} \eta &= \eta_2 - \Lambda_{01} \Lambda_{11}^{-1} \eta_1 \\ \Lambda &= \Lambda_{00} - \Lambda_{01} \Lambda_{11}^{-1} \Lambda_{10} \end{aligned}$$

$$\begin{aligned} \text{Conditioning: } \eta' &= \eta_2 - \Lambda_{01} \Lambda_{11}^{-1} \eta_1 \\ \Lambda' &= \Lambda_{00} \end{aligned}$$

Moment form: Marginal is easy. Condition is hard.

Canonical Form: Marginal is hard. Condition is easy.

You can always transform your ~~information~~ ^{information} in information space and see what happens.

Now a linear transform of Gaussian is still Gaussian.

What about non-linear transform? Of course not Gaussian.

So question comes, How to "preserve" Gaussian for non-linear transform?

2. Uncent Transform.

① Idea: we may not get a Gaussian after the transform. But we can use Gaussian to approximate the transformed distribution.



② How to do this approximation?

- <1> Sample points on G .
- <2> Compute transformed points under g .
- <3> Estimate Gaussian using the transformed points

③ Issues:

- <1> What is a good sampling of G ? Cannot be too much points but must capture important details
- <2> How to sample? What's the algorithm?

④ Criteria: A good sampling.

suppose we have the sampling algorithm. It gives points $x^{[i]}$, $w^{[i]}$ as weight for point $x^{[i]}$. A good sampling would be:

$$\sum_i w^{[i]} = 1$$

$$\mu = \sum_i w^{[i]} x^{[i]}$$

$$\Sigma = \sum_i w^{[i]} (x^{[i]} - \mu)(x^{[i]} - \mu)^T$$

weights should sum to 1 and these samples should be able to reconstruct μ and Σ of G before they transform.
(they capture essence of Gaussian μ, Σ)

⑤ The Uncent transform is all about how to do this sampling.

- <1> points: $x^{[0]} = \mu$. $x^{[i]} = \mu + (\sqrt{(n+1)\Sigma})i$ $i=1 \dots n$
 $x^{[i]} = \mu - (\sqrt{(n+1)\Sigma})i$ $i=n+1 \dots 2n$

Sample $2n+1$ points. First sample center μ . then sample symmetrically. center is where most points are.



What is $(\sqrt{(n+1)\Sigma})i$? i means the i^{th} column.

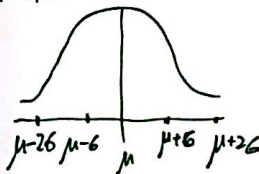
(i) $n+1$ is a scaling factor.

(ii) $\sqrt{\Sigma}$ is square of a matrix

$$\text{General square } \Sigma = U \Lambda U^T = U \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{pmatrix} U^T$$

$$\text{Cholesky decomposition } \Sigma = LL^T. \sqrt{\Sigma} = L.$$

1D intuition



we center around μ . then $\mu \pm \sigma$. then scale factor. before σ .

<2> Weights.

$$w_m^{[0]} = \frac{1}{n+1} \quad w_c^{[0]} = w_m^{[0]} + (1-\alpha)^2 \beta \quad \text{weight for } \mu$$

$$w_m^{[i]} = w_c^{[i]} = \frac{1}{2(n+1)} \quad i=1 \dots 2n. \quad \text{weights for other points}$$

α, β are params one can choose. they control certain things. plot to visualize!

one can plug back and see this satisfies our criterion.

<3> Reconstructed Gaussian (Statistical Estimation)

$$\mu' = \sum_{i=0}^{2n} w_m^{[i]} g(x^{[i]})$$

$$\Sigma' = \sum_{i=0}^{2n} w_c^{[i]} (g(x^{[i]}) - \mu')(g(x^{[i]}) - \mu')^T$$