

PROBABILITY Review

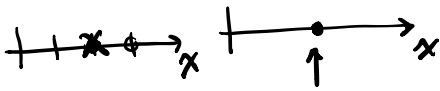
Appendix C | Chapter 2

Random Variables -

1D

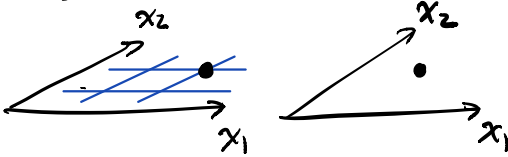
Regular Variable

Discrete Case Continuous Case



2D

Discrete



3D

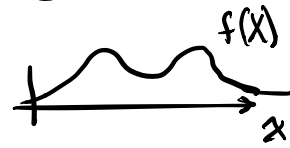
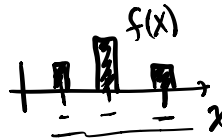


Blowing up points
to density cloud

Random Variable $f(x) \geq 0$

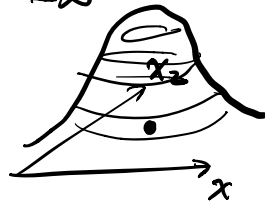
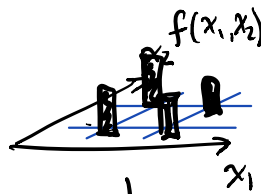
Discrete

Continuous



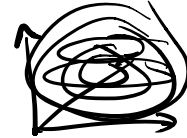
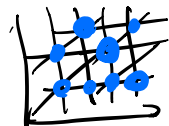
histogram
↓
 $\sum_x f(x) = 1$

and
 $\int_{-\infty}^{\infty} f(x) dx = 1$

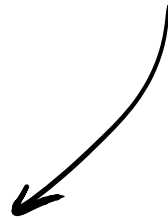


↓
 $\sum_{x_1} \sum_{x_2} f(x_1, x_2) = 1$

$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2) dx_1 dx_2 = 1$

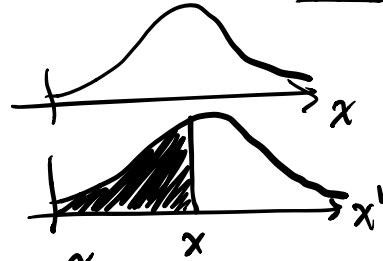


$\sum_{x_1} \sum_{x_2} \sum_{x_3} f(x_1, x_2, x_3) = 1$ $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3) dx_1 dx_2 dx_3 = 1$



How do we define random variables? (continuous)

- density function: $f(x)$
- cumulative density function: $F(x)$



Note: measure theory -
define $F(x)$ first...

$$\frac{\partial F}{\partial x} = f(x)$$

$$\frac{\partial}{\partial x} \rightarrow F(x) = \int_{-\infty}^x f(x') dx'$$

fund. thm of calculus / Leibnitz integration rule

$$\frac{\partial F}{\partial x} = f(x)$$

Expected Value:

"Expected value of $g(x)$ given $f(x)$ "

$$E_{x \sim f}(g(x)) = \int_{-\infty}^{\infty} g(x) f(x) dx$$

Moments: $x \in \mathbb{R}$

" k th moment of f about c "

$$\mu_k = \int_{-\infty}^{\infty} (x-c)^k f(x) dx$$

1D:

- 1st moment of f about 0

$$\mu = \int_{-\infty}^{\infty} x f(x) dx = E_{x \sim f}[x] = E[x]$$

"mean of x "
"expected value of x "

- 2nd moment of f about μ

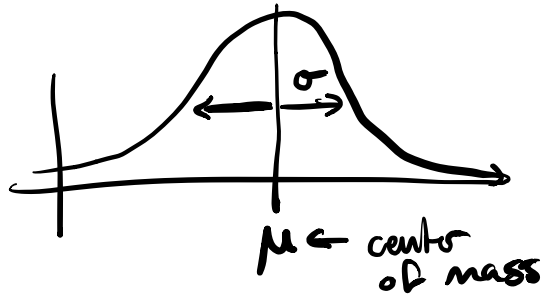
$$\sigma^2 = \int_{-\infty}^{\infty} (x-\mu)^2 f(x) dx$$

"variance of $f(x)$ "

$$\sqrt{\sigma^2} = \sigma = \text{standard deviation}$$

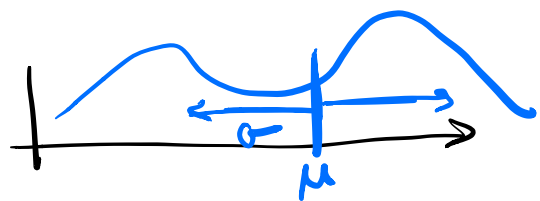
- and so on w higher order terms...

Constructing a Taylor expansion of $f(x)$



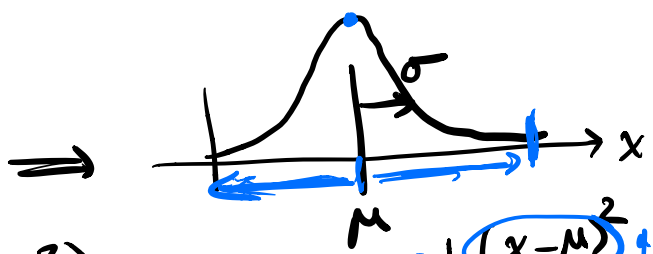
Gaussian distribution:

the mean & variance define $f(x)$
 μ, σ^2 are enough to define $f(x)$ completely



also called

Normal distribution: $x \sim N(\mu, \sigma^2)$



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

2D: $x \in \mathbb{R}^2$ $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ $f: \mathbb{R}^2 \rightarrow \mathbb{R}_+$ $\int f(x) dx = 1$ shape c

Function $g(x)$

Expected of $g(x)$

$$E_{x \sim f}(g(x)) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x) f(x_1, x_2) dx_1 dx_2$$

Moments

- 1st moment about 0
- 2nd moment about μ

$$E(x) = \mu = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f(x_1, x_2) dx_1 dx_2$$

↑
vector

$$E[(x-\mu)(x-\mu)^T] = \sum = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x-\mu)(x-\mu)^T f(x_1, x_2) dx_1 dx_2$$

matrix
covariance matrix

• 3rd moment... (not as important for us...)

3 tensor $M_{ijk} = E((x-\mu)_i (x-\mu)_j (x-\mu)_k) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x-\mu)_i (x-\mu)_j (x-\mu)_k f(x_1, x_2) dx_1 dx_2$

N-Dimensions $x \in \mathbb{R}^n$ $f: x \rightarrow \mathbb{R}_+$

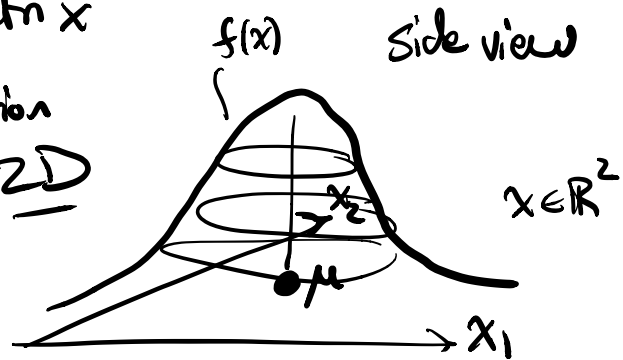
$E(x) = \mu = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} x f(x_1, \dots, x_n) dx_1 \dots dx_n$ ← Mean vector

$E[(x-\mu)(x-\mu)^T] = \Sigma = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (x-\mu)(x-\mu)^T f(x_1, \dots, x_n) dx_1 \dots dx_n$
 matrix covariance matrix

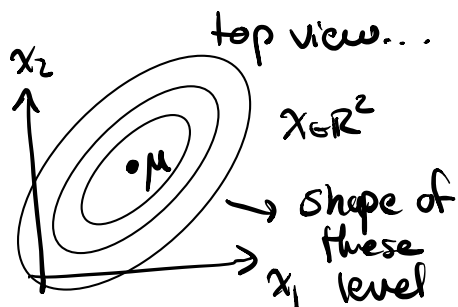
Gaussian (Normal) Distribution

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

2D



$f(x) = \frac{1}{(2\pi)^n \det(\Sigma)} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$



$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$
 shape

Shape of these level set ellipses is given by Σ

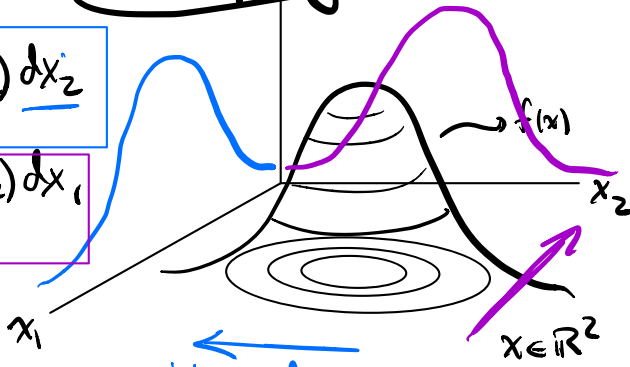
Note:

$f(x_1, \dots, x_n)$: joint distribution

Marginal Distribution \rightarrow collapsing onto a lower dim

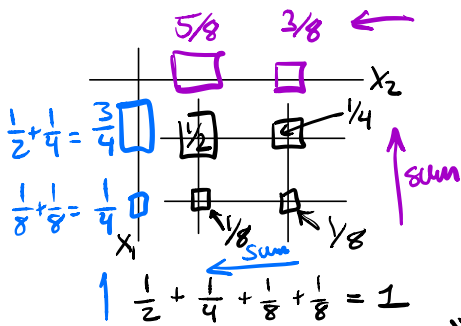
$$f_{x_1}(x) = \int_{-\infty}^{\infty} f(x_1, x_2) dx_2$$

$$f_{x_2}(x_2) = \int_{-\infty}^{\infty} f(x_1, x_2) dx_1$$

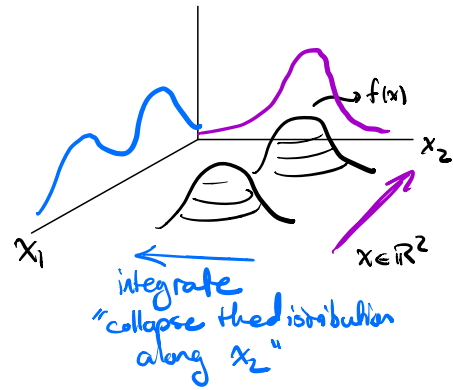


Discrete case
looking down from top.

marginalize \leftarrow "integrate the distribution along x_2 "



"written in the margins."

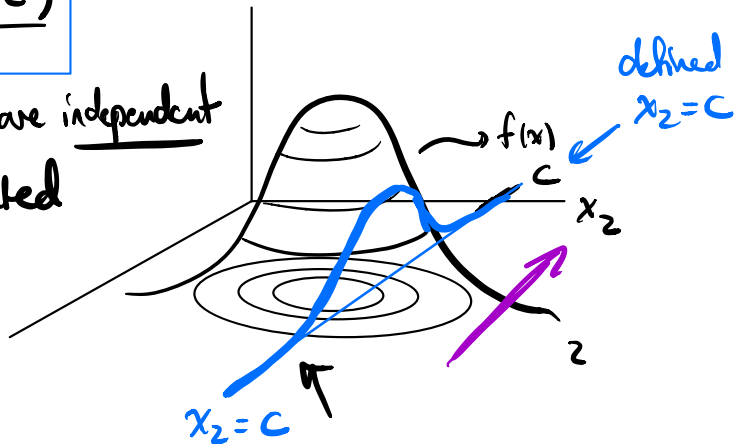


Conditional Distribution Slice along a dimension

$$f(x_1 | x_2 = c) = \frac{1}{\int_{-\infty}^{\infty} f(x_1, c) dx_1} f(x_1, c)$$

$f(x_1, x_2)$: say that x_1 & x_2 are independent
or independently distributed

if
$$f(x_1, x_2) = f_1(x_1) f_2(x_2)$$



Independent Variables x_1, \dots, x_n

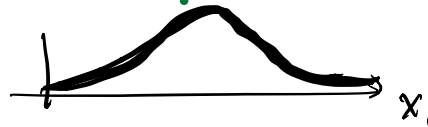
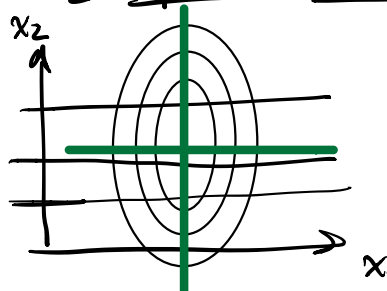
$$f(x_1, \dots, x_n) = \underbrace{f_1(x_1)}_{\text{individual}} \underbrace{f_2(x_2)}_{\text{distributions}} \dots \underbrace{f_n(x_n)}$$

this gives \rightarrow joint distribution 2D:

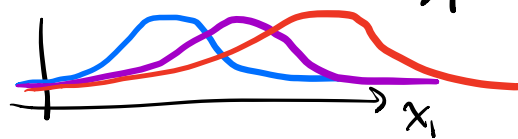
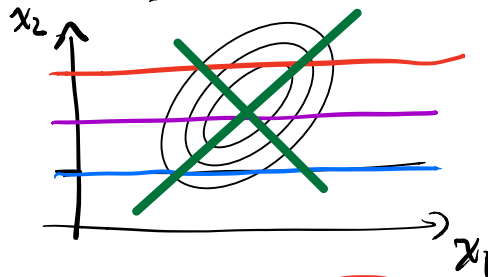
$$f(x_1 | x_2) = \underline{f_1(x_1)}$$

No matter what x_2 is the conditional prob of x_1 doesn't change

joint dist of 2 independent variables



2 dependent variables.



Gaussian level sets

$f(x_1, x_2)$

$$\underline{f(x_1 | x_2=c)} = \frac{1}{\int_{-\infty}^{\infty} f(x_1, c) dx_1} f(x_1, c)$$

$$f(x_1 | x_2=c) = f(x_1, x_2=c)$$

$$\frac{1}{f_2(c) \int_{-\infty}^{\infty} \underbrace{f_1(x_1) f_2(c)}_{=1} dx_1} = \frac{1}{f_2(c)} \underline{f_1(x_1)} \underline{f_2(c)}$$

$$f_{x_1}(x_1) = \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 = f_1(x_1) \underbrace{\int_{-\infty}^{\infty} f_2(x_2) dx_2}_{1} = f_1(x_1)$$

Covariance Matrices $x \in \mathbb{R}^n$

$$\Sigma = E((x - \mu)(x - \mu)^T)$$

$$= E \begin{bmatrix} x_1 - \mu_1 \\ \vdots \\ x_n - \mu_n \end{bmatrix} [x_1 - \mu_1 \dots x_n - \mu_n]$$

$$= E \begin{bmatrix} (x_1 - \mu_1)^2 = \sigma_1^2 & \dots & (x_1 - \mu_1)(x_n - \mu_n) \\ \vdots & \ddots & \vdots \\ (x_n - \mu_n)(x_1 - \mu_1) & \dots & (x_n - \mu_n)^2 = \sigma_n^2 \end{bmatrix}$$

$\rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$
 correlation
 of x_i & x_j

$$\sigma_i^2 = E((x_i - \mu_i)(x_i - \mu_i))$$

$$\sigma_{ij} = E((x_i - \mu_i)(x_j - \mu_j))$$

$\rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$
 correlation
 of x_i & x_j

variance
 terms

$$\sigma_{ij} = \Sigma_{ij} = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (x_i - \mu_i)(x_j - \mu_j) f(x_1, \dots, x_n) dx_1 \dots dx_n$$

what happens if x_1, \dots, x_n are independent?

$$\sigma_{ij} = \sum_{ij} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_i - \mu_i)(x_j - \mu_j) f_i(x_i) f_j(x_j) f_i(x_i) \dots f_n(x_n) dx_i dx_j dx_{\dots}$$

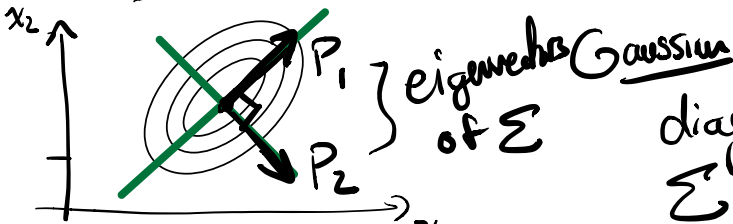
$$\int_{-\infty}^{\infty} (x_i - \mu_i) f_i(x_i) dx_i \int_{-\infty}^{\infty} (x_j - \mu_j) f_j(x_j) dx_j \int_{-\infty}^{\infty} f_i dx_i \int_{-\infty}^{\infty} f_j dx_j \dots$$

if x_i & x_j are independent

$$\Rightarrow \boxed{\sigma_{ij} = \sum_{ij} = 0}$$

if all x_1, \dots, x_n are independent $\Rightarrow \Sigma$ diagonal

2 dependent variables. $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_n^2 \end{bmatrix}$



coord. transform on space P
random vector

$$x = [P_1 P_2] \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

in the exponential \rightarrow orthogonal Σ

$$\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)$$

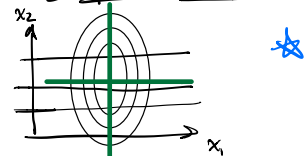
$$f(x) = C e$$

$$\bar{\mu} = P^{-1} \mu$$

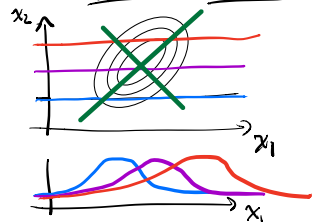
diagonal Σ

non diagonal Σ

2D. joint dist of 2 independent variables



2 dependent variables.



$$\overline{f(z)} \approx e^{-\frac{1}{2}(Pz - P\bar{\mu})^T \Sigma^{-1} (Pz - P\bar{\mu})}$$

$$\overline{f(z)} \approx e^{-\frac{1}{2}(z - \bar{\mu})^T \underbrace{P^T \Sigma^{-1} P}_{\text{diagonalizing } \Sigma} (z - \bar{\mu})}$$

$$\overline{\Sigma}^{-1} = P^T \Sigma^{-1} P \Rightarrow \overline{\Sigma}^{-1} = P \overline{\Sigma}^{-1} P^T$$

diagonal ↓

Σ : sym.

↓
diagonalizable
by an orthonormal
P (rotation)

$$P^T = \begin{bmatrix} p_1^T \\ p_2^T \end{bmatrix}$$

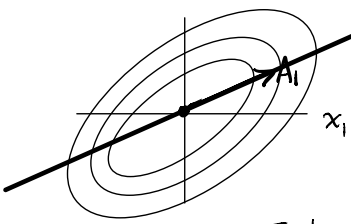
↑
left eigen
vectors

$$[p_1, p_2]$$

right eigenvectors
of Σ

Σ & Σ^{-1}
have the
same eigenvectors

Slicing a Gaussian Distribution 2×2



$$x \sim \mathcal{N}(0, \Sigma)$$

$$x = A_1 z \rightarrow$$

\uparrow
 \mathbb{R}^2

$$f(x) \approx e^{-\frac{1}{2}(x-0)^T \Sigma^{-1}(x-0)} \quad z \sim \mathcal{N}(0, \Sigma_z)$$

\downarrow
 1×1

$$e^{-\frac{1}{2} z A_1^T \Sigma^{-1} A_1 z}$$

$\underbrace{\hspace{10em}}_{\Sigma_z^{-1}}$

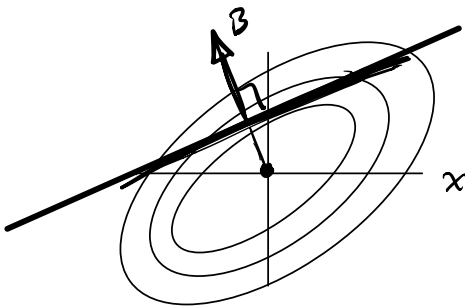
still Gaussian in z

$$\Sigma_z^{-1} = A_1^T \Sigma^{-1} A_1$$

\leftarrow $\left[\begin{array}{c} \parallel \\ \parallel \end{array} \right]$

$$\Sigma_z = (A_1^T \Sigma^{-1} A_1)^{-1}$$

Constraints on x .



$$B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \in \mathbb{R}^2$$

$$B^T x = 1$$

$$\begin{bmatrix} B_1 & B_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 1$$

$$B_1 x_1 + B_2 x_2 = 1$$

Multivariate Gaussians
are "jointly gaussian"

$$x_1 = \frac{1 - B_2 x_2}{B_1}$$

"Any slice of a jointly Gaussian distribution is still jointly Gaussian"

\Rightarrow lin combs of joint Gauss rand. variables are joint. Gauss.

$$\approx e^{-\frac{1}{2} (x^T \Sigma^{-1} x)}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1/B_1 \\ 0 \end{bmatrix} + \begin{bmatrix} -B_2/B_1 \\ 1 \end{bmatrix} x_2$$

$$e^{-\frac{1}{2} \left(\begin{bmatrix} 1/B_1 \\ 0 \end{bmatrix} + \begin{bmatrix} -B_2/B_1 \\ 1 \end{bmatrix} x_2 \right)^T \Sigma^{-1} \left(\begin{bmatrix} 1/B_1 \\ 0 \end{bmatrix} + \begin{bmatrix} -B_2/B_1 \\ 1 \end{bmatrix} x_2 \right)}$$

$\underbrace{\hspace{10em}}_{x_2}$

Matrix Derivatives:

$$f(x) \quad x \in \mathbb{R}^{m \times n} \quad \frac{\partial f}{\partial x} = ?$$

before $\frac{\partial f}{\partial x}$
↑
vector

1. Vectorize X. \Rightarrow "stack X into a vector"
(stack up the cols of X)

$$\text{vec}(X) \in \mathbb{R}^{mn}$$

$$\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B) \quad \leftarrow \text{wikipedia}$$

$\frac{\partial \text{vec} f}{\partial \text{vec} X}$ ↑ \otimes ← Kronecker product

Ex. $f(Q) = x^T Q x$ \Rightarrow $\frac{\partial \text{vec} f}{\partial \text{vec} Q} = x^T \otimes x^T$

Scalar

$$\text{vec}(f) = \text{vec}(x^T Q x) = (x^T \otimes x^T) \text{vec} Q$$

A ↑ C ↑

Ex. $f(x) = Ax$ $\frac{\partial \text{vec} f}{\partial \text{vec} x} = I \otimes A$

matrix

$$\text{vec}(f) = \text{vec}(AxI) = (I \otimes A) \text{vec}(x)$$

f

2. think about derivative differently.
doing derivative elementwise...

Recall

$$f(x) = C^T x$$

\uparrow \uparrow \uparrow
 row col vector

$$\frac{\partial f}{\partial x} = C^T \quad \Delta f = \frac{\partial f}{\partial x} \Delta x$$

General version:

what is $\langle \cdot, \cdot \rangle$
for matrices?

$$\Delta f = \left\langle \frac{\partial f}{\partial x}, \Delta x \right\rangle$$

define derivative
to do this...

$$\langle A, B \rangle = \sum_{ij} A_{ij} B_{ij}$$

$$\langle A, B \rangle = \sum_{ij} A_{ij} B_{ij} = \underline{\text{Tr}(A^T B)}$$

Many times:

$$A = [A_1 \dots A_n]$$

$$\begin{bmatrix} A_1^T \\ \vdots \\ A_n^T \end{bmatrix} [B_1 \dots B_n]$$

$$f(x) = \text{Tr}(X)$$

$$B = [B_1 \dots B_n]$$

$$A_1^T B_1 \dots A_n^T B_n$$

$$\text{Tr}(ABC) = \text{Tr}(CAB) = \text{Tr}(BCA)$$

Ex. $f(x) = \text{Tr}(AXB) = \text{Tr}(BAX)$

$$= \langle A^T B^T, X \rangle$$

$$\frac{\partial f}{\partial x} = A^T B^T$$

Useful Relationships:

element wise: $\frac{\partial f}{\partial x_{ij}} = [A^T B^T]_{ij}$

$$\frac{\partial \text{Tr}(BAC)}{\partial A} = B^T C^T$$

$$\text{Tr}(BAC) = \text{Tr}(CBA) = \langle \underline{B^T C^T}, A \rangle$$

$$\frac{\partial \text{Tr}(ABA^T)}{\partial A} = A(B+B^T)$$

page 68 of book

$$f(A) = \text{Tr}(ABA^T) \quad \Delta f = \langle \underline{\quad}, \underline{\Delta A} \rangle$$

$$\Delta f = \text{Tr}(\underline{\Delta A} B A^T) + \text{Tr}(A B \underline{\Delta A}^T)$$

$$\text{Tr}(M) = \text{Tr}(M^T)$$

$$\downarrow \text{Tr}(\underline{\Delta A} B^T A^T)$$

$$= \text{Tr}(\Delta A (B+B^T) A^T)$$

$$= \text{Tr}(\underline{(B+B^T) A^T} \Delta A) = \langle \underline{A(B+B^T)}, \Delta A \rangle$$

$$\frac{\partial f}{\partial A}$$

BACK TO ESTIMATION Chapter 2

PROBABILITY PERSPECTIVE ON LS. ↙

LINEAR MEAS. $\tilde{y} = Hx + v$ $v \sim N(0, R)$

FIND LINEAR ESTIMATOR

$\hat{x} = M\tilde{y} + n$

unbiased

- unbiased (★)
- minimum variance (★)

$E[\hat{x}] = E[x]$
est. \uparrow true

if not $E[\hat{x} - x] = \text{BIAS.}$

$$\begin{aligned} E[\hat{x}] &= E[MHx + Mv + n] \\ &= E[MHx] + E[Mv] + E[n] \\ &= MHE[x] + ME[v] + \bar{n} \end{aligned}$$

$E[\hat{x}] = MHE[x] + n$

} linearly
of $E[\cdot]$
(from integral)

\Rightarrow $MH = I$, $n = 0$

$$\hat{x} = My \quad \dots \quad \swarrow \text{sum of variance terms}$$

$$\min_M J = \frac{1}{2} E \left[\frac{(\hat{x} - x)(\hat{x} - x)^T}{\sum_i (\hat{x}_i - x_i)^2} \right] = \frac{1}{2} \text{Tr} E \left[\underbrace{(\hat{x} - x)(\hat{x} - x)^T}_{\text{covariance of } \hat{x} - x} \right]$$

s.t. $MH = I \quad (\hat{x} = My)$

$$\downarrow$$

$$\min_M J = \frac{1}{2} \text{Tr} E \left[(\hat{x} - x)(\hat{x} - x)^T \right]$$

s.t. $MH = I$

Parallel Axis Theorem: (for unbiased estimator)

"how does covariance shift when you shift the mean of a distribution"

axis of rotation \swarrow moment of inertia

unbiased $\hat{x} = \frac{MHx}{I}$

$$E \left[(\hat{x} - x)(\hat{x} - x)^T \right] = E \left[\hat{x}\hat{x}^T \right] - E \left[x\hat{x}^T \right] - E \left[\hat{x}x^T \right] + E \left[xx^T \right]$$

with respect to distribution of \hat{x}

$$= E \left[\hat{x}\hat{x}^T \right] - E \left[xx^T \right] \frac{H^T H}{I} - \frac{MH}{I} E \left[xx^T \right] + E \left[xx^T \right]$$

$$= E \left[\hat{x}\hat{x}^T \right] - E \left[xx^T \right]$$

Minimize variance

$$\min_M J = \frac{1}{2} \text{Tr}(\underline{E[\hat{x}\hat{x}^T]} - \underline{E[xx^T]})$$

$$\text{s.t. } \underline{MH = I} \Rightarrow (\underline{I - MH}) = 0$$

Now OPTIMIZE ...

$$\mathcal{L}(M, \Lambda) = \frac{1}{2} \text{Tr}(\underline{E[\hat{x}\hat{x}^T]} - \underline{E[xx^T]}) + \langle \underline{\Lambda^T}, \underline{I - MH} \rangle$$

↑
primal

↑
dual

$$\frac{\partial \mathcal{L}}{\partial M} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \Lambda} = 0$$

$$\text{Tr}(\Lambda(I - MH))$$

$$\hat{x} = My = M(Hx + \underline{v}) = \underline{MH}x + \underline{Mv}$$

$$E[\hat{x}\hat{x}^T] = E(MHxx^TH^TM^T) + E(MHxv^TM^T) + E(Mvx^TH^TM^T)$$

$$MH = I \quad + E(Mvv^TM^T)$$

$$= \underline{E[xx^T]} + \underline{E[xv^T]}M^T + M\underline{E[vx^T]}$$

$$+ M\underline{E[vv^T]}M^T$$

x, v are independent

noise to be independent from parameters (x)

$$MRM^T$$

$$\underline{E(xv^T)} = \underline{E(x)E(v^T)^0} \quad \star \quad v \sim N(0, R)$$

$$\mathcal{L}(M, \Lambda) = \frac{1}{2} \text{Tr}(MRM^T) + \text{Tr}(\Lambda(\mathbf{I} - MH))$$

$$\frac{\partial \mathcal{L}}{\partial M} = \frac{1}{2} \frac{\partial}{\partial M} (\text{Tr}(MRM^T) - \text{Tr}(\Lambda MH))$$

$$\frac{\partial}{\partial A} \text{Tr}(BAC) = B^T C^T$$

$$\frac{\partial}{\partial A} \text{Tr}(ABA^T) = A(B+B^T)$$

$$= \frac{1}{2} M(R+R^T) - \Lambda^T H^T$$

$$\rightarrow \text{Tr}(\Lambda MH)$$

$$= MR - \Lambda^T H^T = 0$$

$$\Rightarrow M = \Lambda^T H^T R^{-1} \leftarrow \text{LQR term...}$$

$$\frac{\partial \mathcal{L}}{\partial \Lambda} = \frac{\partial}{\partial \Lambda} \text{Tr}(\Lambda(\mathbf{I} - MH))$$

$$= (\mathbf{I} - MH)^T = 0 \Rightarrow MH = \mathbf{I}$$

Combining ...

$$(\Lambda^T H^T R^{-1}) H = \mathbf{I} \Rightarrow \Lambda^T = (H^T R^{-1} H)^{-1}$$

$$M = (H^T R^{-1} H)^{-1} H^T R^{-1}$$

$$\hat{x} = M\tilde{y} = (H^T R^{-1} H)^{-1} H^T R^{-1} \tilde{y}$$

← minimum variance linear estimator

Note: has the same form as weighted least squares

where $W = R^{-1}$

weight matrix = inverse of the covariance of noise

if you have noise $v \sim N(0, R)$

$$\tilde{y} = Hx + v$$

optimal way to weight LS is to invert covariance R ... use R^{-1} as a weighting matrix.

Next time: what if you have a prior estimate on x

$$x = \hat{x}_a + w \quad w \sim N(0, Q)$$