

Example Distributions

| Distribution | PDF or PMF | Mean | Variance |
|---------------------------|--|---------------------|-----------------------|
| $Bernoulli(p)$ | $\begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0. \end{cases}$ | p | $p(1 - p)$ |
| $Binomial(n, p)$ | $\binom{n}{k} p^k (1 - p)^{n-k} \text{ for } k = 0, 1, \dots, n$ | np | $np(1 - p)$ |
| $Geometric(p)$ | $p(1 - p)^{k-1} \text{ for } k = 1, 2, \dots$ | $\frac{1}{p}$ | $\frac{1-p}{p^2}$ |
| $Poisson(\lambda)$ | $\frac{e^{-\lambda} \lambda^k}{k!} \text{ for } k = 0, 1, \dots$ | λ | λ |
| $Uniform(a, b)$ | $\frac{1}{b-a} \text{ for all } x \in (a, b)$ | $\frac{a+b}{2}$ | $\frac{(b-a)^2}{12}$ |
| $Gaussian(\mu, \sigma^2)$ | $\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \text{ for all } x \in (-\infty, \infty)$ | μ | σ^2 |
| $Exponential(\lambda)$ | $\lambda e^{-\lambda x} \text{ for all } x \geq 0, \lambda \geq 0$ | $\frac{1}{\lambda}$ | $\frac{1}{\lambda^2}$ |

Two R.V.

Bivariate CDF

$$F_{XY}(x, y) = P(X \leq x, Y \leq y)$$

Bivariate PMF

$$P_{XY}(x, y) = P(X=x, Y=y)$$

Marginal

Given $\rightarrow P(x, y) \rightarrow$ Captures max info

$$\iint p(x,y) dx dy = 1$$

$$p(x) = \sum_y p(x,y) = \int p(x,y) dy$$

$$p(y) = \sum_x p(x,y) = \int p(x,y) dx$$

distribution wrt x

$$p(x,y) = p(x) p(y|x)$$

$$= p(\text{Marginal}) \quad p(\text{Conditional})$$

$$p(x)p(y|x) = p(y)p(x|y)$$

$$p(y|x) = \frac{p(y)p(x|y)}{p(x)}$$

$$= \frac{p(y)p(x|y)}{\sum_{y'} p(y') p(x|y')}$$

Two r.v. are independent if:

$$P_{XY}(x, y) = P_X(x) P_Y(y)$$

$$P_{Y|X}(x, y) = P_Y(y)$$

$$\begin{aligned}\text{Cov}[X, Y] &= E[(X - E(X))(Y - E(Y))] \\ &= E[XY - XE(Y) - YE(X) + E(X)E(Y)] \\ &= E(XY) - \cancel{E(X)E(Y)} - \cancel{E(Y)E(X)} + \cancel{E(X)E(Y)} \\ &= E(XY) - E(X)E(Y)\end{aligned}$$

$$\begin{aligned}\text{Var}[X+Y] &= \text{Var}[X] + \text{Var}[Y] \\ &\quad + 2\text{Cov}[X, Y]\end{aligned}$$

$$\text{Var}(X) = E[X^2] - E[X]^2$$

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

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

$$\xrightarrow{\text{normalizing constant}} (2\pi)^{n/2} |\Sigma|^{\frac{1}{2}}$$

$x \rightarrow$ space over which it is defined

$\mu, \Sigma \rightarrow$ parameters
 \uparrow
 cov matrix

Σ is always P.S.D.

Proof: $X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad P(x)$

$$\Sigma = \int_{\mathbb{R}^n} P(x) (x - \mu)(x - \mu)^T dx$$

$$\mu = E(x) = \int_{\mathbb{R}^n} P(x) x dx$$

$$\text{then } v^T \Sigma v \geq 0 \text{ as } \stackrel{?}{=}$$

If Gaussian n dimension

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

Multivariant Gaussian (Normal examples)

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

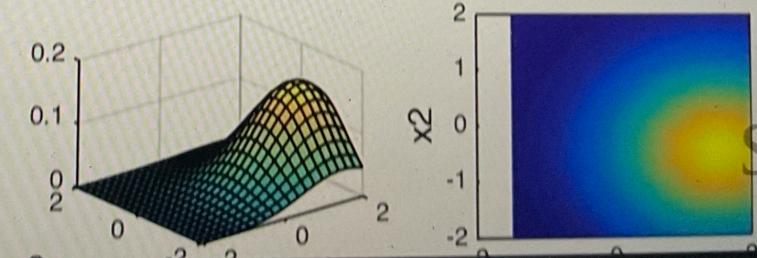
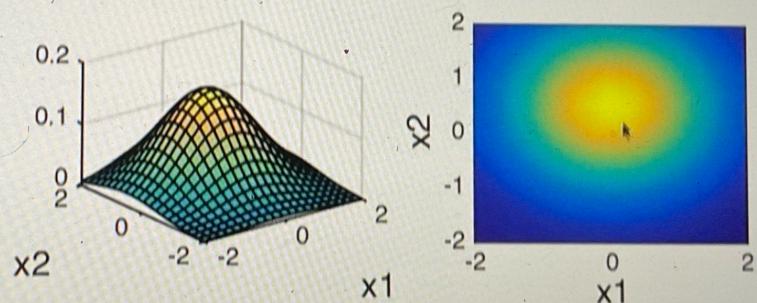
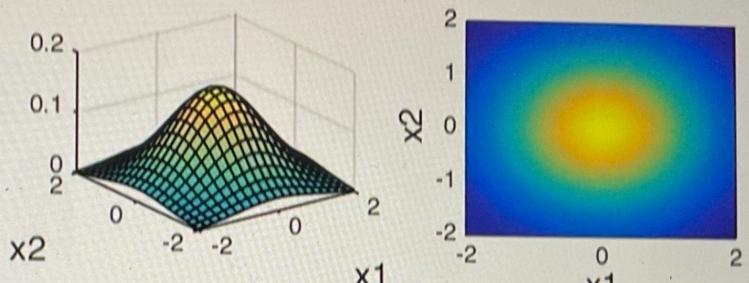
$$\mu = [0 \ 0]^T$$

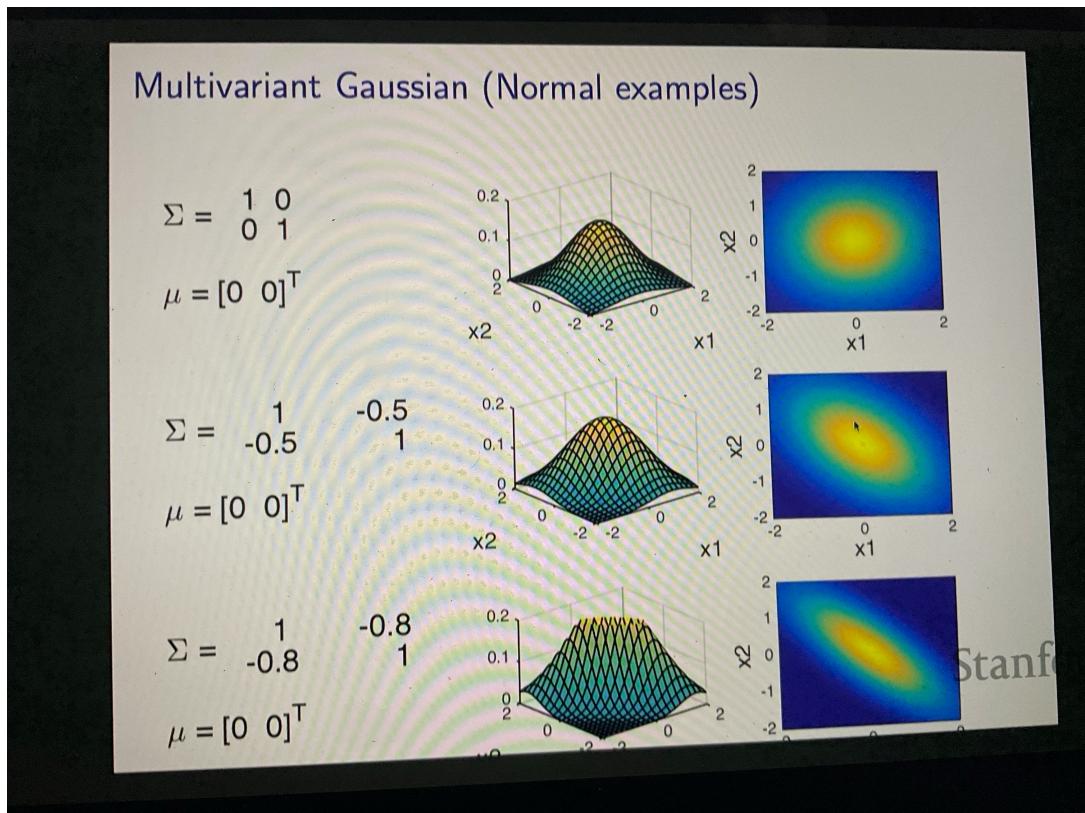
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0 \quad 0.5]^T$$

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [1.5 \quad -0.5]^T$$



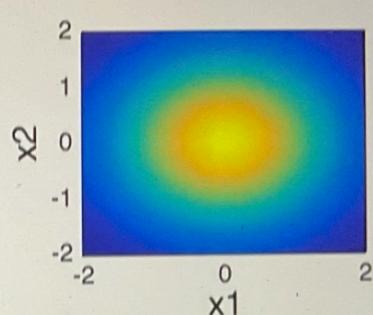
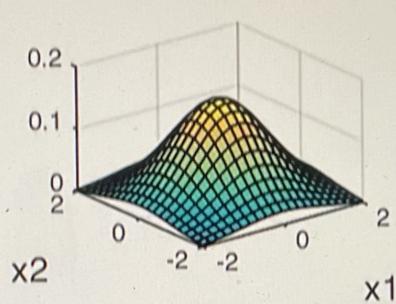


Stanf

Multivariant Gaussian (Normal examples)

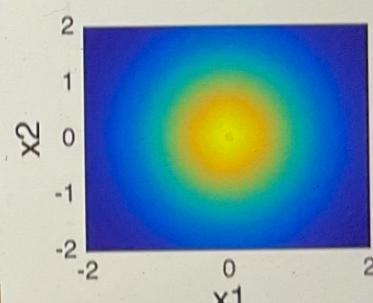
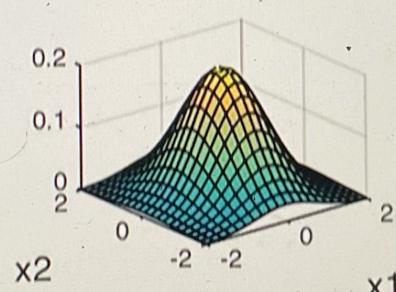
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0 \ 0]^T$$



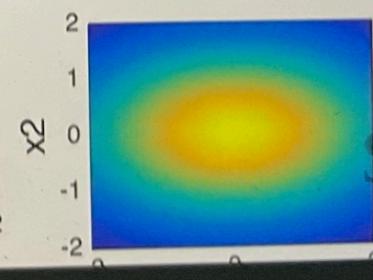
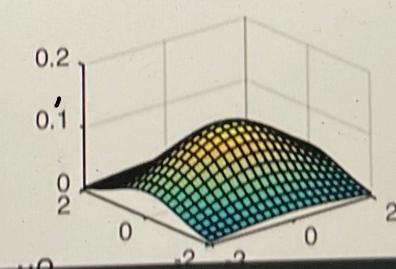
$$\Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0 \ 0]^T$$

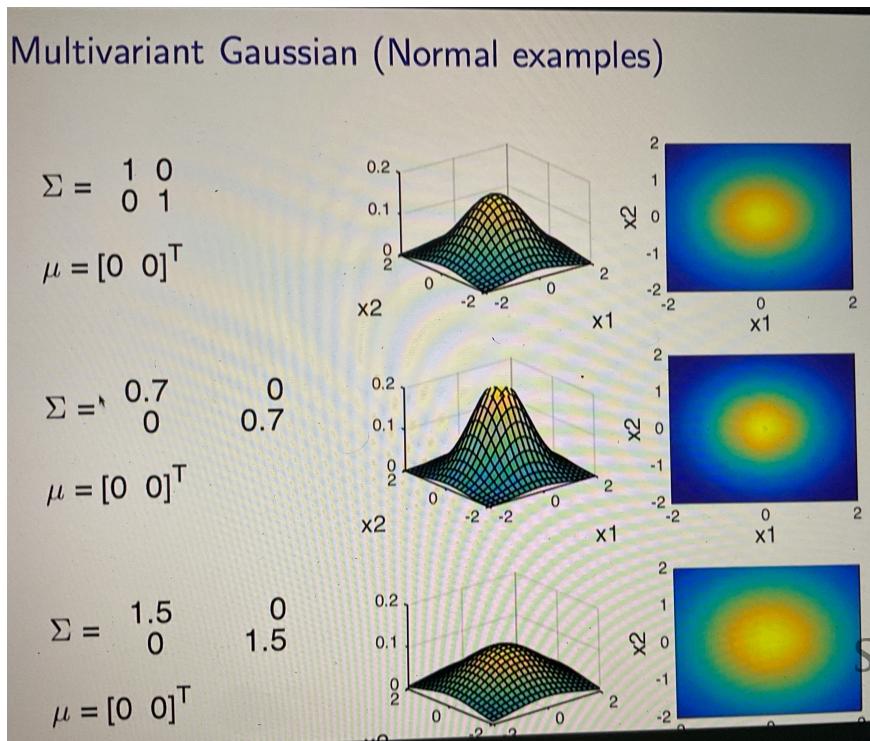


$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0 \ 0]^T$$



Stanf



$$P(X=x|Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

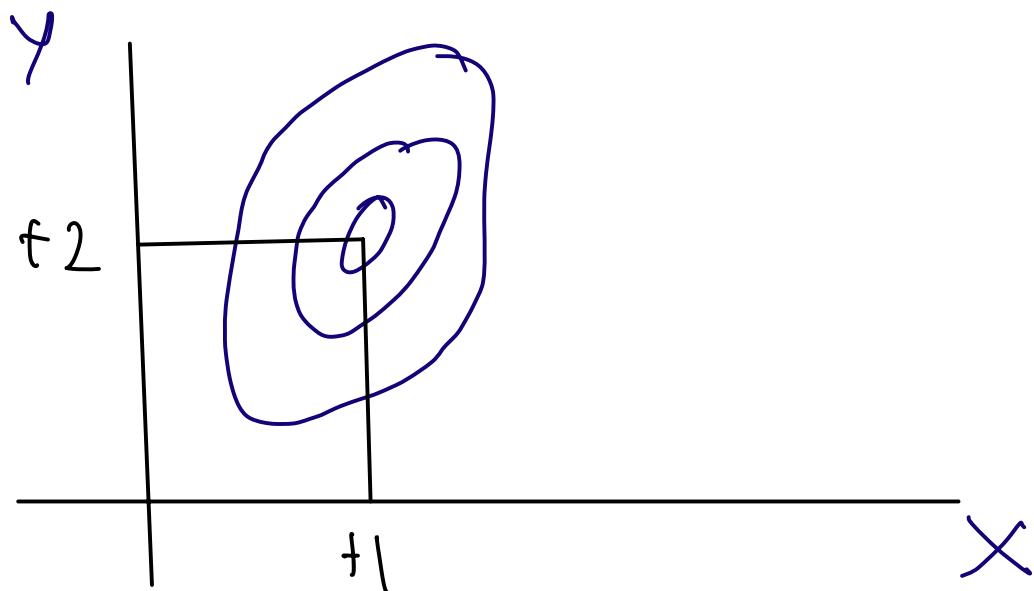
$$E(X|Y=y) = \sum_x \frac{x P(X=x, Y=y)}{P(Y=y)}$$

$E[X|Y](y) = E[X|Y=y]$ is a func of y
 so it's R.V.

$E[X]$ - constant

$E[X|Y]$ - R.V. over Y space (sample space of Y)

$E[X|Y=y]$ → function of y



$E[X|Y=y]$

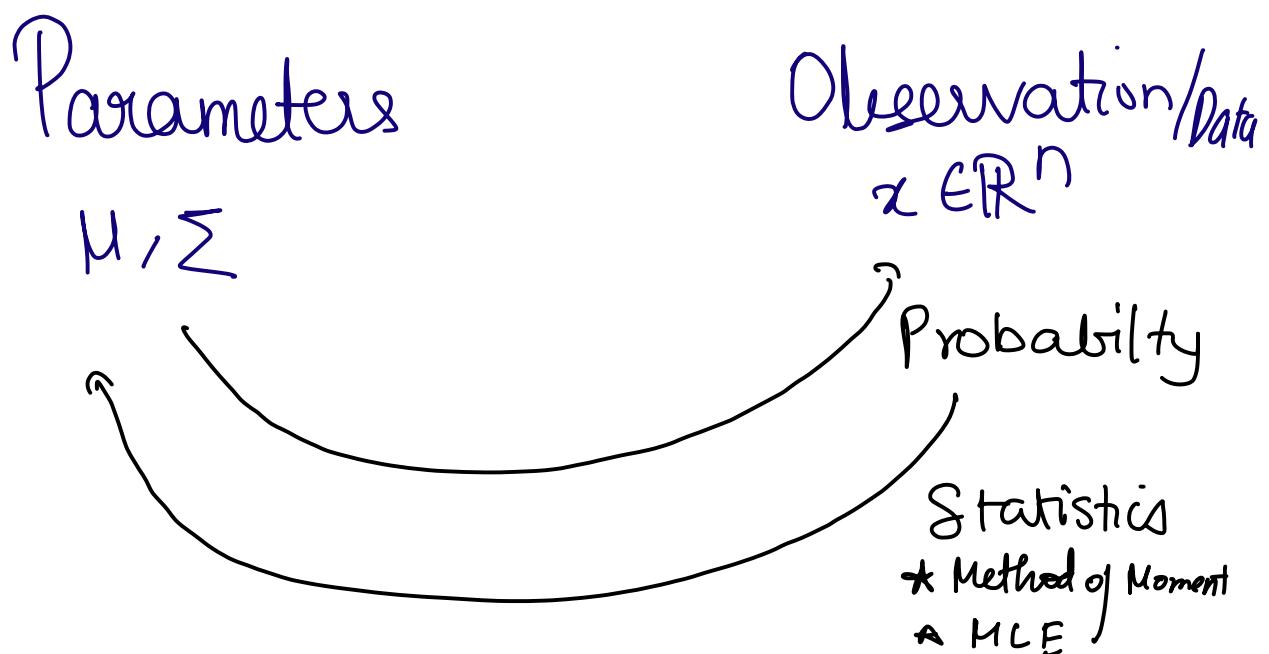
LOTUS

$$E[X] = E[E[X|Y]]$$

$$\begin{aligned}
E[E(X|Y)] &= E \left[\sum_x x P(X=x|Y) \right] \\
&= \sum_y \left(\sum_x x P(X=x|Y=y) \right) P(Y=y) \\
&= \sum_y \sum_x x P(X=x, Y=y) \\
&= X \sum_y I(X=x, Y=y) \\
&= \sum_x x P(X=x) \\
&= E[X]
\end{aligned}$$

$$P(a|b, c) = \frac{P(b|a, c)P(a|c)}{P(b|c)}$$

$$P(a|b) = \frac{P(b|a) P(a)}{P(b)}$$



Training Data

$(x, y) \rightarrow$ learn parameters (Stat)

Model → predict future data
(Prob)

GAUSSIAN

$x \in \mathbb{R}^d$

$X = \{x^{(1)}, \dots, x^{(n)}\}$ $x^{(i)} \rightarrow$ i^{th} example
i.i.d. $\in \mathbb{R}^d$

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

$$p(x^{(1)}, \dots, x^{(n)}) = \prod_{i=1}^n p(x^{(i)})$$

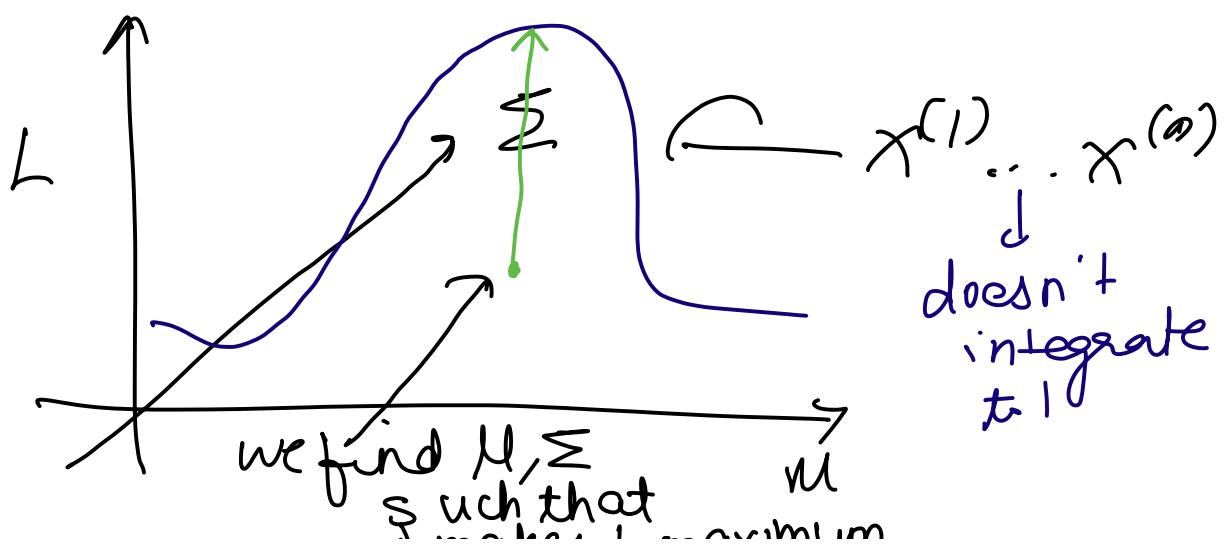
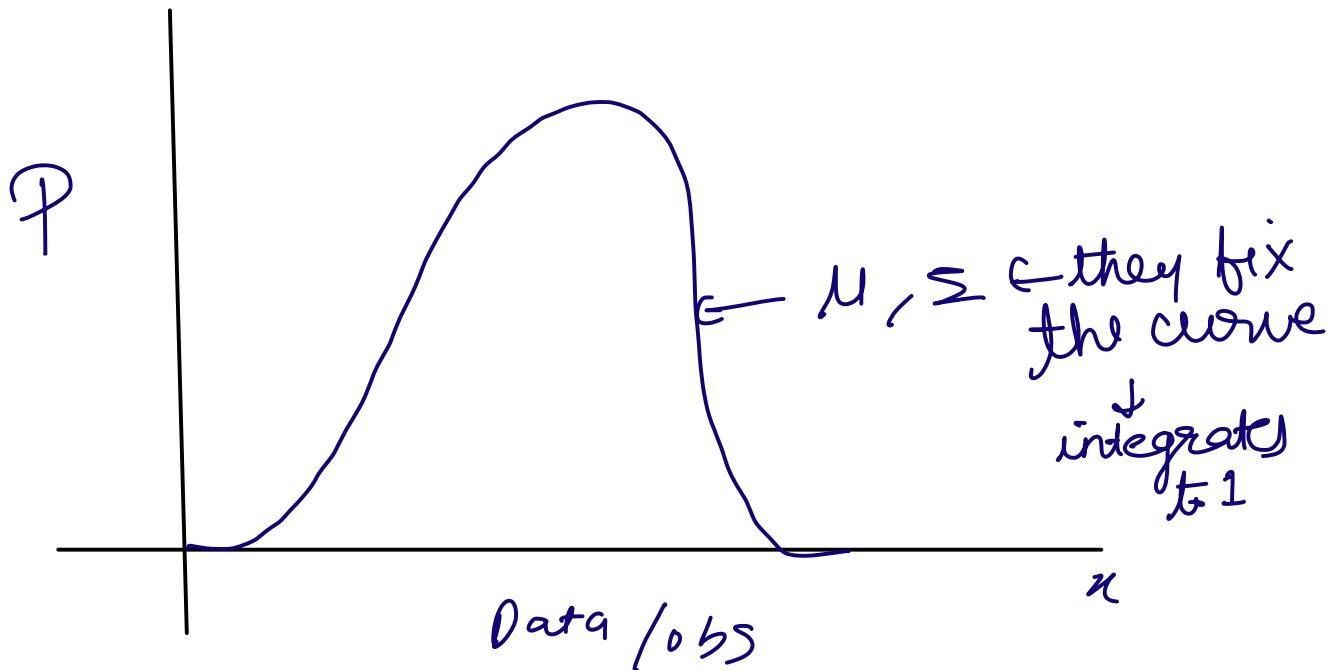
$$p(x^{(1)}, \dots, x^{(n)}; \mu, \Sigma) = \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)\right)$$

likelihood

$$L(\mu, \Sigma; x^{(1)}, \dots, x^{(n)}) =$$

[Same expression, different interpretation)

probability of data given parameters
likelihood of parameters given data



it makes L maximum

$$L(\theta; X) = \prod_{i=1}^n L(\theta; x^{(i)})$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \prod_{i=1}^n L(\theta; x^{(i)})$$

$$= \arg \max_{\theta} \log \prod_{i=1}^n L(\theta; x^{(i)})$$

$$= \arg \max_{\theta} \sum_{i=1}^n \log L(\theta; x^{(i)})$$

$$= \arg \max_{\theta} \sum_{i=1}^n \ell(\theta; x^{(i)})$$

GAUSSIAN

$$= \arg \max_{\mu, \Sigma} \sum_{i=1}^n \log \left[\frac{1}{(2\pi)^{d/2}} |\Sigma|^{1/2} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right) \right]$$

$$\hat{\mu}, \hat{\Sigma} = \sum_{i=1}^n \mathbf{x} - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}^{(i)} - \mu)^T \Sigma^{-1} (\mathbf{x}^{(i)} - \mu)$$

For μ

$$\nabla_{\mu} \left(\sum_{i=1}^n k - 0.5 \log |\Sigma| - 0.5 (\mathbf{x}^{(i)} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}) \right)$$
$$= \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu}) - 0.5 \times 2 \times (\mathbf{x}^{(i)} - \boldsymbol{\mu}) \Sigma^{-1}$$

.

$$\sum_{i=1}^n (\Sigma^{-1} \mathbf{x}^{(i)} - \Sigma^{-1} \boldsymbol{\mu}) = 0$$

$$n \Sigma^{-1} \boldsymbol{\mu} = \sum_{i=1}^n \Sigma^{-1} \mathbf{x}^{(i)}$$

$$\Sigma^{-1} \boldsymbol{\mu} = \sum^{-1} \sum_{i=1}^n \frac{1}{n} \mathbf{x}^{(i)}$$

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$

For Σ

$$S = \sum^{-1} \quad \nabla_S l = 0 \Leftrightarrow \nabla_{\sum^{-1}} l = 0$$

$$\begin{aligned} & \nabla_S \left(\sum_{i=1}^n \frac{1}{2} \log |S| - \frac{1}{2} (x^{(i)} - \mu)^T S (x^{(i)} - \mu) \right) \\ &= \frac{1}{2} \left[n S^{-1} - \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T \right] \\ \sum &= S^{-1} = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T \\ \text{Use } & \nabla_A x^T A x = x x^T \end{aligned}$$

Proof:

$$\langle A, B \rangle = \text{tr}(A^T B) = A : B$$

$$\begin{aligned} \langle A, BC \rangle &= A : BC = ACT : B \\ &= B^T A : C \\ &= C^T B^T : A^T \\ &= BC : A \end{aligned}$$

$$d(x^T A x)$$

$$d(x : Ax)$$

$$d(x x^T : A)$$

$$\tilde{X}X^T : dA'$$

$$\Sigma \approx E(X - E(X))$$

$$d(X^T A X)$$

$$d(X : A X)$$

$$dX : A X + X : A dX$$

$$dX : A X + A^T X : dX$$

$$dX : \underline{\underline{A X + A^T X}}$$

$$\text{Note: } L(\theta) = L(\theta; X, \vec{y}) = P(\vec{y} | X; \theta)$$