

Bayesian Learning

- ① What is a probability distribution?
- ② What is a random variable?
- ③ Different types of probability distributions

— Bernoulli, Binomial, Poisson

— Normal / Gaussian, Beta,

For any probability distribution

— What is its domain? X

— What is its pmf or pdf.

— What are the parameters.

$$E[x] = \sum_{x \in X} P(x=x)$$

or

$$\int_x x p(x) dx$$

$$E[f(x)] = \sum_{x \in X} f(x) P(x=x)$$

or

$$\int f(x) p(x) dx$$

Why deal with a probability distribution?

Poisson Distribution

$$X = \{0, 1, 2, \dots\}$$

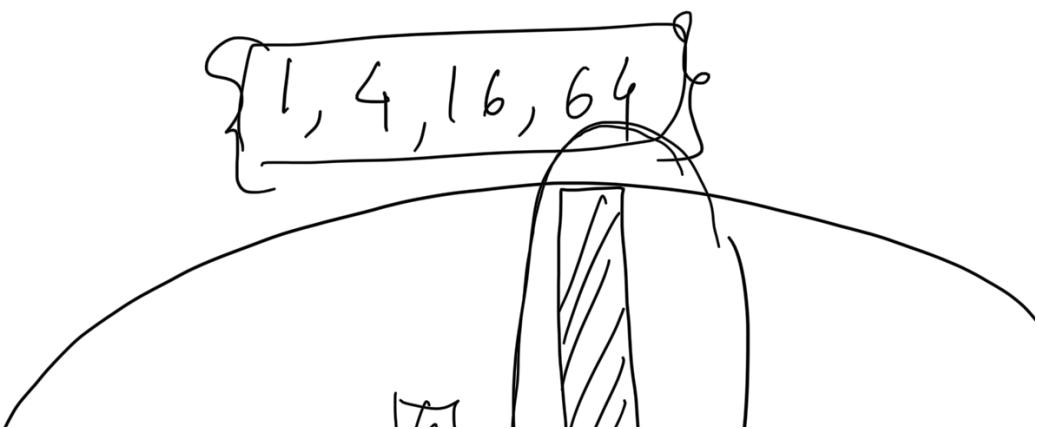
$$P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

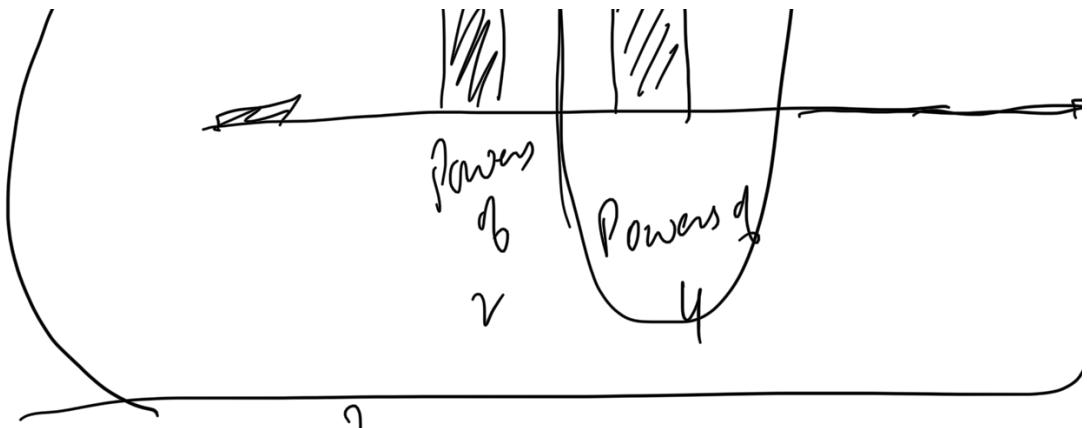
$$\lambda = 4$$

$$E[X] = \lambda$$

What are the parameters of the distribution?

$$\{1, 4, 16, 64\}$$





Why?

(all numbers between 1 & 100)

Concept = all numbers between 1 & 100
1, 4, 16, 64

Concept = all powers of 4

choose concept with highest
likelihood

$$D = \{ \underline{1, 4, 16, 64} \}$$

$h = \text{All powers of 4}$

$$P(1/h) = \frac{1}{4}$$

$$D(4/h) = 1$$

$$P(16|h) = \frac{1}{4}$$

$$P(64|h) = \frac{1}{4}$$

$$L(D|h) = P(x_1|h) P(x_2|h) \dots P(x_N|h)$$

likelihood \rightarrow
 $= \prod_{x \in D} P(x|h)$

$$\log(L(D|h)) = \sum_{x \in D} \log P(x|h)$$

$$L(D|h) = \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} .$$

Let $h = \text{all even numbers}$,

$$D = \{1, 4, 16, 64\}$$

$$P(1|h) = \frac{1}{50}$$

$$L(D|h) = \left(\frac{1}{50}\right)^4$$

Let $h = \text{multiples of } 10$

- 10 20 30 40 50

P(1 | m)

Calculate the likelihoods of D

① given every h .

Pick h with highest likelihood.

Maximum likelihood Estimation
(MLE)

② Use some prior information

and pick h with highest prior.

Max a-Priori estimation

③ Combine likelihood and the prior.

Maximum a-Posteriori Estimation
(MAP)

$P(h) \rightarrow$ Prior prob that h is
the hypothesis.

$P(D|h) \rightarrow$ likelihood

$$P(h|D) =$$

↑
posterior

$$\frac{P(D|h)}{\sum_{h' \in H} P(D|h') p(h')}$$

$$\frac{P(D|h)}{\sum_{h' \in H} P(D|h') p(h')}$$

normalization

$$\sum_{h \in \mathcal{H}} p(h' | D) = 1$$

normal-
constant

h^* - is the right
hypothesis

h prior \rightarrow All even numbers

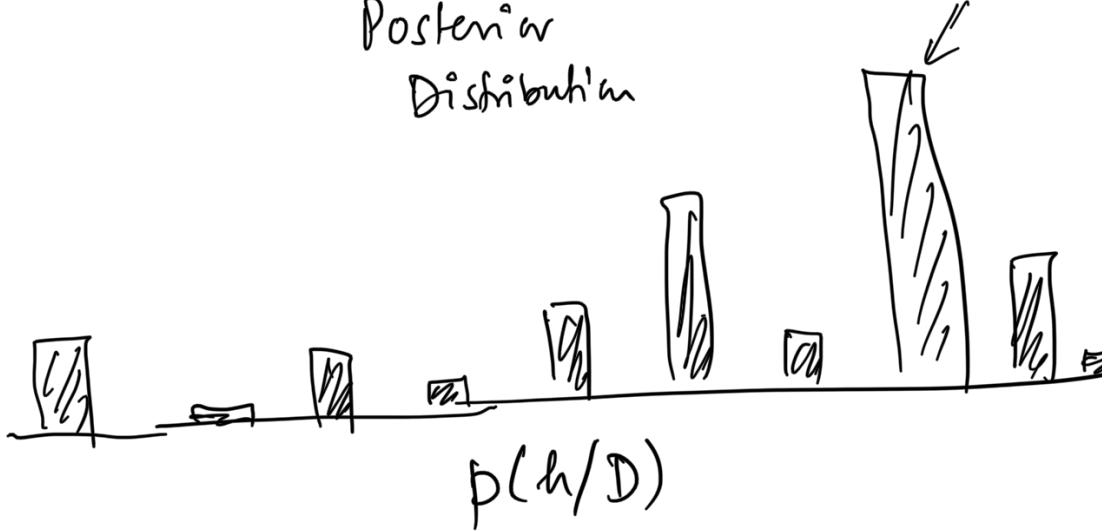
$$p(x^* = 12 | D) = p(x^* = 12 | \text{h prior})$$



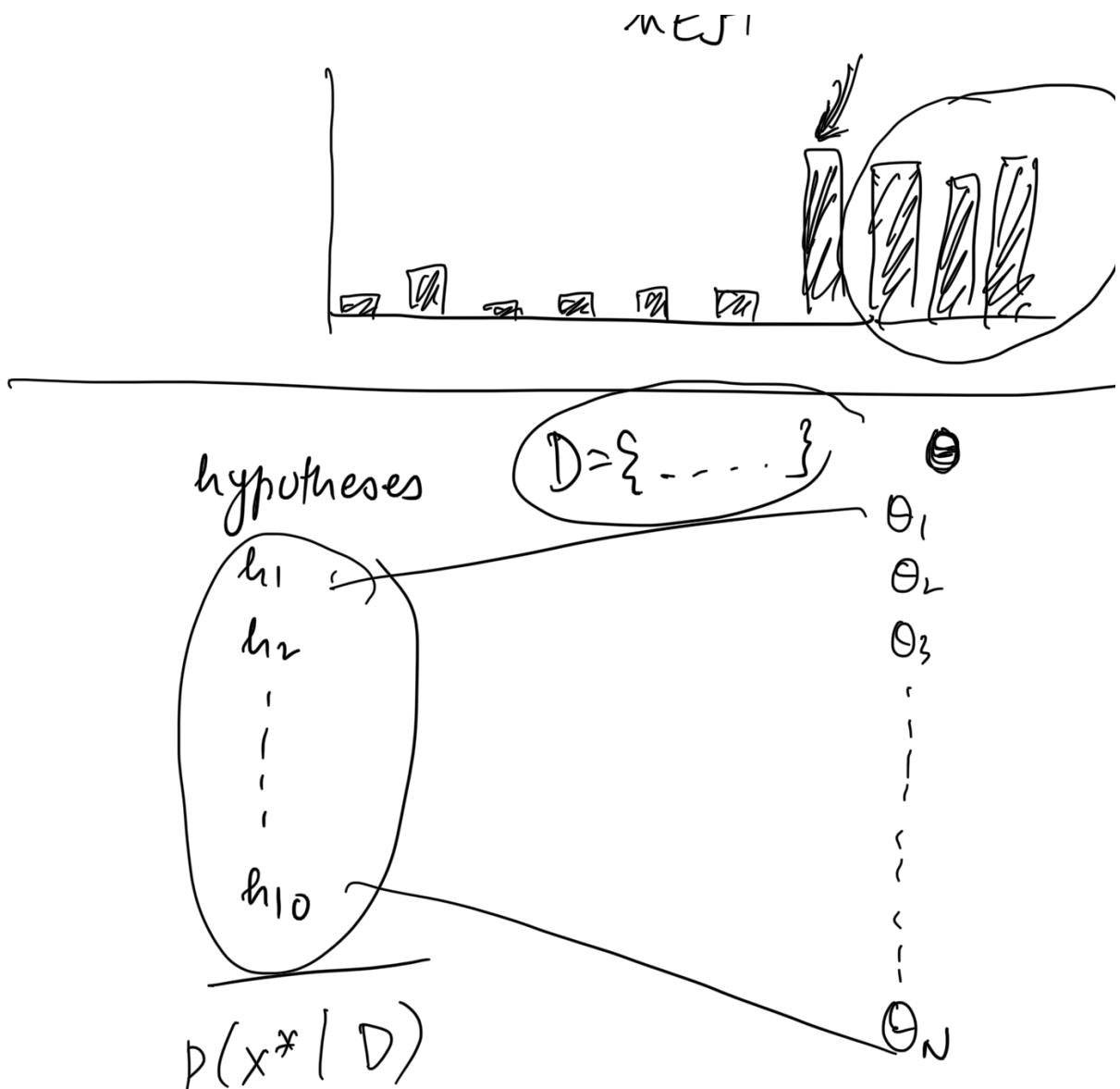
=

Bayesian Averaging

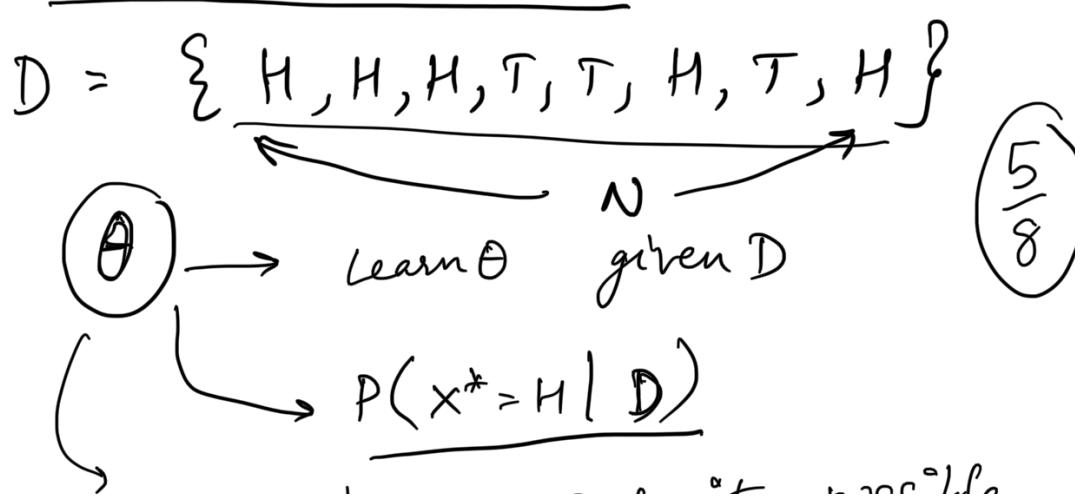
Posterior
Distribution



$$p(x^* = 12 | D) = \sum_{h \in \mathcal{H}} p(x^* = 12 | h) p(h | D)$$



Bernoulli Distribution



$0 \leq \theta \leq 1 \rightarrow$ infinite possible values.

Likelihood of D given θ .

$$\text{likelihood} = \prod_{x \in D} P(x=x | \theta)$$

$$\left. \begin{aligned} P(x=H | \theta) &= \theta \\ P(x=T | \theta) &= 1 - \theta \end{aligned} \right]$$

Let D have N_1 heads & N_0 tails

$$N_1 + N_0 = N$$

$$L(D | \theta) = \prod P(x=x | \theta)$$

$$= \frac{\theta^{N_1} (1-\theta)^{N_0}}{N!}$$

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} L(D | \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \frac{\theta^{N_1} (1-\theta)^{N_0}}{N!}$$

$$\hat{\theta}_{MLE} = \frac{N_1}{N_0 + N_1} = \frac{N_1}{N}$$

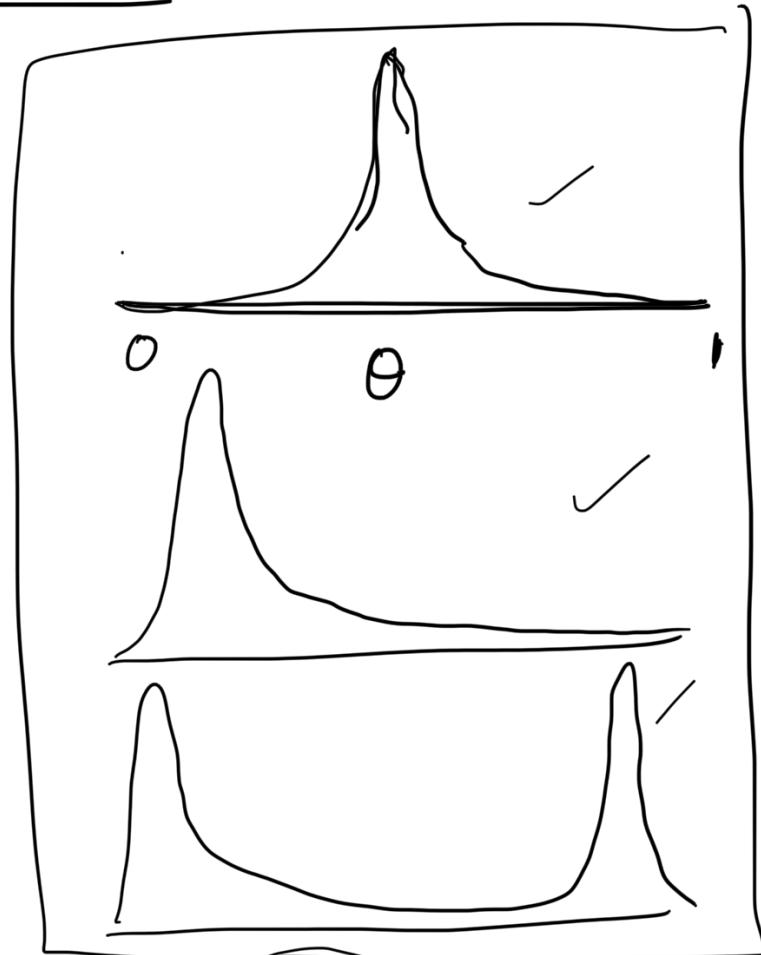
$$D / x^* = H | n) = P(x^* = H | \hat{\theta}_{MLE})$$

$$\hat{\theta}_{MLE} = \frac{N_1}{N}$$

Incorporating Prior

Beta

a, b



Posterior =

~~Likelihood~~ \times Prior

$$P(D|\theta) = \theta^{N_1} (1-\theta)^{N_0}$$

$$P(\theta) \propto \theta^a (1-\theta)^{b-1}$$

$$\begin{aligned}
 p(\theta | D) &= \frac{\theta^{N_1} (1-\theta)^{N_0} \theta^{a-1} (1-\theta)^{b-1}}{\int_{\theta'}^1 \theta'^{N_1} (1-\theta')^{N_0} \theta'^{a-1} (1-\theta')^{b-1} d\theta'} \\
 &= \frac{\theta^{N_1+a-1} (1-\theta)^{N_0+b-1}}{\int_0^1 d\theta'}
 \end{aligned}$$

$p(\theta | D)$ is also a Beta distribution
with parameters are $\boxed{N_1+a, N_0+b}$

If prior & likelihood pdf are conjugate pairs

Then the posterior will have the same form as the prior.

Now we know that the posterior

$$p(\theta | D) = \text{Beta}(\theta | a+N_1, b+N_0)$$

$$\hat{\theta}_{MAP} = \frac{a+N_1-1}{a+b+N-2}$$

$N = N_0 + N_1$

$\hat{\theta}$

$a-1$

$$\hat{\theta}_{\text{Prior}} = \frac{a+b-2}{a+b-2}$$

$$\hat{\theta}_{\text{MLE}} = \frac{N_1}{N_0+N_1}$$

$$\hat{\theta}_{\text{MAP}} = \frac{a+N_1-1}{a+b+N-2} \quad N = N_0+N_1$$

$$P(X^*=H|D) \rightarrow \hat{\theta}_{\text{Prior}}$$

$$P(X^*=H|D) \rightarrow \hat{\theta}_{\text{MLE}}$$

$$P(X^*=H|D) \rightarrow \hat{\theta}_{\text{MAP}}$$

$$\int P(X^*=H|\theta) P(\theta|D) d\theta$$

$$= \frac{a+N_1}{a+b+N} \quad N = N_0+N_1$$

Black Swan Paradox

$$D = \{T, T, T\}$$

$$N_1 = 0 \quad N_0 = 3 \quad N_1 + N_0 = 3$$

$$\hat{\theta}_{\text{MLE}} = 0$$

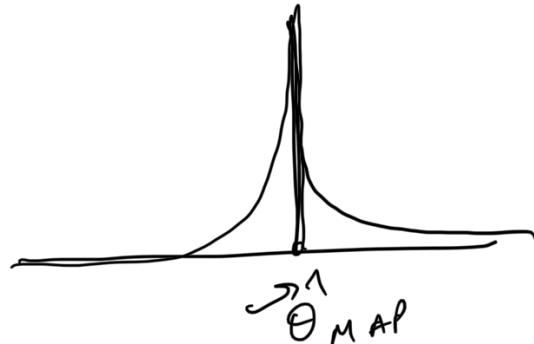
$$P(X^*=H|\hat{\theta}_{\text{MLE}}) = 0$$

... n - 1

$$P(x_{\text{H}}^* | \hat{\theta}_{\text{MAP}}) = \frac{a+n_1}{a+b+N-2}$$

$a = 2$ $b = 10$

$$P(x_{\text{H}}^* | \hat{\theta}_{\text{MAP}}) = \frac{2+0-1}{2+10+3-2} = \frac{1}{13}$$



$$\text{Var}[\theta | D] = \frac{(a+n_1)(b+n_0)}{(a+b+N)^2(a+b+N+1)}$$

MVN

μ $D \times 1$ vector

Σ $D \times D$ matrix

$$D = \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \rightarrow D \times 1 \text{ vector}$$

$$\begin{aligned}
 & \overbrace{\vdots}^x_N \\
 \text{pdf}(x | \mu, \Sigma) &= \frac{1}{(2\pi)^N |\Sigma|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] \\
 L(D | \mu, \Sigma) &= \prod_{i=1}^N \text{pdf}(x_i | \mu, \Sigma) \\
 LL(D | \mu, \Sigma) &= \sum_{i=1}^N \log \text{pdf}(x_i | \mu, \Sigma) \\
 \boxed{\begin{aligned}
 \hat{\mu}_{MLE} &= \frac{1}{N} \sum_{i=1}^N x_i \triangleq \bar{x} \\
 \hat{\Sigma}_{MLE} &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T
 \end{aligned}}
 \end{aligned}$$

- Choosing a Gaussian ~~prior~~ prior on μ will give a Gaussian posterior for μ (I-Wishart) distribution

① MAP for θ

② MLE for univariate Gaussian

③ MLE for multivariate Gaussian.

$$X - \{H, T\}$$

$$\theta \quad 0 \leq \theta \leq 1$$

$$P(X=H) = \theta$$

$$P(X=T) = 1-\theta$$

Prior on $\theta \rightarrow$ Beta distribution

$$\theta \quad 0 \leq \theta \leq 1$$

$$pdf(\theta) = \text{Beta}(\theta | a, b)$$

$$= \frac{1}{\text{Beta}(a, b)} \theta^{a-1} (1-\theta)^{b-1}$$

$$\mathbb{E}[\theta] = \int \theta p(\theta) d\theta = \int_0^1 \theta \frac{1}{\text{Beta}(a, b)} \theta^{a-1} (1-\theta)^{b-1}$$
$$= \frac{a}{a+b}$$

What θ gives the max pdf

Mode

$$\arg \max_{\theta} p_{\text{af}}(\theta)$$

$$\arg \max_{\theta} \theta^{a-1} (1-\theta)^{b-1}$$

$$\begin{aligned} \frac{d}{d\theta} \theta^{a-1} (1-\theta)^{b-1} &= (a-1) \theta^{a-2} (1-\theta)^{b-1} \\ &\quad + \theta^{a-1} (b-1) (1-\theta)^{b-2} \\ &= (a-1) \theta^{a-2} (1-\theta)^{b-1} \\ &\quad - (b-1) \theta^{a-1} (1-\theta)^{b-2} \end{aligned}$$

Setting to 0 & solving for θ

$$(a-1) \theta^{a-2} (1-\theta)^{b-1} = (b-1) \theta^{a-1} (1-\theta)$$

$$(a-1)(1-\theta) = (b-1)\theta$$

$$\theta [b-1 + a-1] = a-1$$

$$\boxed{\hat{\theta} = \frac{a-1}{a+b-2}}$$

Posterior for θ : $p(\theta | D)$

$$= \text{Beta}(\theta | a+N_1, b+N_0)$$

$\uparrow N_1 - \# \text{heads in } D$
 $N_0 - \# \text{tails in } D$

$$\hat{\theta}_{MAP} = \frac{a+N_1-1}{a+N_1+b+N_0-2} = \frac{a+N_1-1}{a+b+a-1}$$

univariate

bivariate

$$N = N_1 + N_0$$

Univariate Gaussian Distribution

$$X \quad -\infty \leq x \leq \infty$$

$$paf(x) = N(x | \mu, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right]$$

$$D = \{x_1, x_2, x_3, \dots, x_N\}$$

i.i.d independent

& identically distributed

$$L(D) = \prod_{i=1}^N p(x_i)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x_i-\mu)^2\right]$$

$$\ell \ell(D) = \log L(D)$$

$$= \sum_{i=1}^N \log []$$

$$= \sum_{i=1}^N \left[-\frac{1}{2} \log(2\pi) - \log \sigma - \frac{1}{2\sigma^2} (x_i-\mu)^2 \right]$$

$$\begin{aligned}
 &= -\frac{N}{2} \log(\frac{1}{2}\sigma^2) - N \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 \\
 \frac{\partial}{\partial \mu} \underline{\underline{\ell \ell(D)}} &= -\frac{1}{\sigma^2} \sum_{i=1}^N 2(x_i - \mu)(-1) \\
 &= \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) \\
 &= \frac{1}{\sigma^2} \sum_{i=1}^N x_i - \frac{1}{\sigma^2} N\mu
 \end{aligned}$$

Set to 0 and solve for μ

$$\boxed{\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i} \quad \text{Sample mean or average}$$

$$\frac{\partial}{\partial \sigma} \underline{\underline{\ell \ell(D)}} = -\frac{N}{\sigma} - \frac{1}{2\sigma^3} \sum_{i=1}^N (x_i - \mu)^2$$

Set to 0 and solve for σ

$$\boxed{\hat{\sigma}^2_{MLE} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2} \quad \text{Sample variance}$$

$$D = \{x_1, x_2, x_3, \dots, x_N\}$$

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

$$p.d.f(x_i) = \mathcal{N}(x_i | \mu, \Sigma)$$

$$\rightarrow \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} [(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)] \right\}$$

$$\ell \ell(\theta) = -\frac{Nd}{2} \log(2\pi) - \frac{N}{2} \log |\Sigma|$$

$$= \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\frac{\partial}{\partial \mu} \ell \ell(\theta) = -\frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \mu} [(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)]$$

$$\frac{d}{da} a^T A a = (A + A^T) a$$

$$= -\frac{1}{2} \sum_{i=1}^N (\Sigma^{-1} + \Sigma^{-T}) (x_i - \mu) (-1)$$

$$\frac{\partial}{\partial \mu} \ell \ell(\theta) = 0$$

$$\sum_{i=1}^N (x_i - \mu) \geq 0$$

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\ell l(D) = -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\Lambda = \Sigma^{-1}$$

$$|\Lambda| = \frac{1}{|\Sigma|}$$

— cook book

$$\ell l(D) = \frac{N}{2} \log |\Lambda| - \frac{1}{2} \sum (x_i - \mu)^T \Lambda (x_i - \mu)$$

"Trace Trick"

$$a^T A a = \underbrace{\text{tr}(a a^T A)}_{d \times 1 \times 1 \times d \quad d \times d}$$

$$\text{tr}(B) = \sum_{i=1}^d B_{ii}$$

$$\begin{aligned} \ell l(D) &= \frac{N}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^N \text{tr}[(x_i - \mu)(x_i - \mu)^T \Lambda] \\ &= \frac{N}{2} \log |\Lambda| - \frac{1}{2} \text{tr}\left[\underbrace{\left[\sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T \right]}_{S_\mu} \right] \Lambda \end{aligned}$$

$$\underbrace{\sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T}_{= S_\mu} = \text{scatter matrix}$$

$$\ell l(D) = \frac{N}{2} \log |\Lambda| - \frac{1}{2} \text{tr}(S_\mu \Lambda)$$

$$\frac{d \ell(D)}{d \Lambda} = \frac{N}{2} \frac{d}{d \Lambda} \log |\Lambda| - \frac{1}{2} \frac{d}{d \Lambda} \text{tr}(S_{\mu} \Lambda)$$

$$\frac{d}{d \Lambda} \log |\Lambda| = (\Lambda^{-1})^T$$

$$\frac{d}{d \Lambda} \text{tr}(B \Lambda) = B$$

$$\frac{d \ell(D)}{d \Lambda} = \frac{N}{2} (\Lambda^{-1})^T - \frac{1}{2} S_{\mu}$$

$$\Lambda^{-1} = \Sigma \quad \Sigma^T = \Sigma$$

$$\frac{d \ell(D)}{d \Lambda} = \frac{N}{2} \Sigma - \frac{1}{2} S_{\mu}$$

Set it to 0, solve for Σ

$$\begin{aligned} \Sigma &= \frac{1}{N} S_{\mu} \\ &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T \end{aligned}$$

↑
sample covariance matrix