How do we find a prob. dist for a r.v. X?

3 steps

(1) choose a parametric model (eg Gaussian)

$$\theta = \text{parameters}$$

(2) collect samples from r.v. X:

$$D = \{ x_1, \ldots, x_N \}$$

we assume $x_i$ are i.i.d. samples

(3) MLE (maximum likelihood principle)

The optimal parameter $\theta^*$ is that which maximises the probability (likelihood) of the training data.

$$\theta^* = \text{argmax}_\theta \; p(D|\theta)$$

"likelihood function"

Likelihood of data w.r.t. parameter $\theta$

Note ① D is known, so $p(D|\theta)$ is a function of $\theta$

It is not a prob. w.r.t $\theta$

② log = ln

$$= \text{argmax}_\theta \; \log p(D|\theta)$$

$\ell(\theta) = $ log. likelihood. function.

$$= -\text{argmin}_\theta \; -\log p(D|\theta)$$

negative log likelihood function (loss)

data LL

$$\ell(\theta) = \log p(D|\theta)$$

$$= \sum_{i=1}^{N} \log P_i(D_i|\theta)$$ 

} assume independence (i.i.d.)

To get optimal MLE Solution:

○ if $\theta$ is a scalar, at local optimal
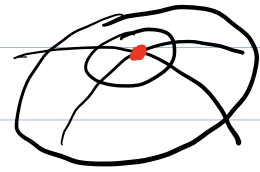
1) $\frac{\partial}{\partial \theta} \log p(D|\theta) = 0$ at $\theta^*$

2) $\frac{\partial^2}{\partial \theta^2} (\log p(D|\theta)) < 0$ at $\theta^*$ (concave)

3)    check the boundary condition of $\theta$ (if necessary)

○   if $\theta$ is a vector

   i) $\nabla_\theta \, l(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} l(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} l(\theta) \end{bmatrix} = 0$

   ii) Hessian Matrix

$$\nabla^2_\theta \, l(\theta) = \begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} & \cdots & \frac{\partial^2}{\partial \theta_1 \partial_p} \\ & \ddots & \\ & & \frac{\partial^2}{\partial \theta_p^2} \end{bmatrix} l(\theta) \quad \begin{cases} 0 \; \begin{pmatrix} negative \\ definite \end{pmatrix} \end{cases}$$

$H \prec 0$   negative definite :   $\theta^T H \theta < 0 \, , \forall \theta$

          $\Rightarrow$ for all directions, "Concave" (mountain)

$H \succ 0$   positive definite :   $\theta^T H \theta > 0 \, , \forall \theta$

                     "convex" (bowl)

( semi-negative-definite)     "ridge"


Ex.   Bernoulli

$$\theta = \pi \in \{0, 1\}$$

$$l(\theta) = \sum_{i=1}^{N} \log P(x_i | \theta)$$

$$= \sum_{i=1}^{N} \log \left( \pi^{x_i} (1-\pi)^{(1-x_i)} \right)$$

$$= \sum_{i=1}^{N} x_i \log \pi + (1-x_i) \log (1-\pi)$$

$$= \left( \sum_{i=1}^{N} x_i \right) \log \pi + \left( \sum_{i=1}^{N} (1-x_i) \right) \log (1-\pi)$$

         # of 1s                # of 0s

$$m = \sum_{i=1}^{N} x_i \; \leftarrow \text{ "sufficient statistics"}$$

                     (a set of statistics) you need

$$= m \log \pi + (1-m) \log (1-\pi)$$

1) find max

$$\frac{\partial}{\partial \pi} \ell(\theta) = \frac{m}{\pi} + \frac{N-m}{1-\pi}(-1) = 0$$

$$\Rightarrow \quad \pi = \frac{m}{N} \quad \hat{\pi} = \frac{m}{N} \text{ (sample mean)}$$

2) $$\frac{\partial^2}{\partial \pi^2} \ell(\theta) = -\frac{m}{\pi^2} - \frac{N-m}{(1-\pi)^2} < 0$$

3) boundary condition: $0 \le m \le N$   $0 \le \frac{m}{N} \le 1$

· Ex. Gaussian

① $\theta = \mu$ ($\alpha^2$ known)

$$\ell(\theta) = \sum_i \left( -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \alpha^2 - \frac{1}{2\alpha^2}(x_i - \mu)^2 \right)$$

$$= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \alpha^2 - \frac{1}{2\alpha^2} \sum_i (x_i - \mu)^2$$

"sufficient statistics" $\Rightarrow$

$$\left\{ \sum_i x_i, \sum_i x_i^2 \right\}$$

i) $$\frac{\partial \ell(\theta)}{\partial \mu} = -\frac{1}{2\alpha^2} \sum_i 2(x_i - \mu)(-1) = 0$$

$$\Rightarrow \quad \hat{\mu} = \frac{\sum_{i=1}^{N} x_i}{N}$$

② $\theta = \alpha^2$

$$\frac{\partial \ell(\theta)}{\partial \alpha^2} = -\frac{N}{2} \frac{1}{\alpha^2} + \frac{1}{\alpha^2}(\ldots) \sum_i (x_i - \mu)^2 = 0$$

$$2a^? \quad - \quad -\frac{}{2}a^2 \quad - \quad 2a^4 \; (-1)\sum(x_i - \mu) = 0$$

$$\hat{\sigma^2} = \frac{1}{N}(x_i - \mu)^2 \;(\text{sample variance})$$

---

Evaluation

Estimate (e.g. $\hat{\mu}$, $\hat{\sigma}_2$) is a number

Estimator is a r.v. (over possible datasets)
$$f(x_1, \ldots, x_N) = \frac{1}{N} \sum_{i=1}^{N} x_i$$

<span style="color:green">(MLE is one of the ways → computing estimators)</span>

<span style="color:red">↑ r.v. for each sample
$x_i \sim p(x_i | \theta)$ true distribution
(assume $x_i$ is drawn from some distribution)</span>

The estimate is the value of the estimator
for a given dataset $D$
$$\hat{\mu} = f(x_1, \ldots, x_N)\big|_{x_1 = x_1, \ldots} \;\; : \;\; \frac{1}{N} \sum_i x_i$$
<span style="color:red">↑ sample</span>

Since the estimator is a r.v., we can derive
the mean & variance to qualify "goodness"
Bias & Variance $\quad \hat{\theta} = f(x_1, \ldots, x_N)$

1) will it converge to the true value of $\theta$?
$$\text{Bias}(\hat{\theta}) = E_{x_1 \ldots x_N}[\hat{\theta} - \theta] = \underline{E_x[\hat{\theta}]} - \theta$$
$\qquad\qquad\qquad\qquad\qquad \uparrow$
$\qquad\qquad\qquad\qquad \text{true value} \quad \text{mean of estimator}$

<span style="color:green">if the bias is non-zero, then we can never get the
true value (even if infinite samples)</span>

2) How long will it take to converge

(How many samples do we need)

$$Var(\hat{\theta}) = E_{x_1 \cdots x_N}[(\hat{\theta} - E\hat{\theta})^2]$$

Ex. Gaussian.

Estimator $\hat{\mu} = \frac{1}{N} \sum_i x_i$

Mean of $\hat{\mu}$   $E_{x_1 \cdots x_N}[\frac{1}{N} \sum_i x_i] = \frac{1}{N} N \cdot \mu = \mu$

Bias$(\hat{\mu}) = 0$

Var of $\hat{\mu}$   $E_{x_1 \cdots x_N}[(\hat{\mu} - E\hat{\mu})^2]$

$$= E_{x_1 \cdots x_N}[(\frac{1}{N} \sum_i x_i - E\hat{\mu})^2]$$

$(a+b)^2 = a^2 + ab + ba + b^2$

$$= \frac{1}{N^2} E((\sum_i (x_i - \mu))^2)$$

$(\sum_i x_i)^2 = \sum_i x_i \sum_j x_j$

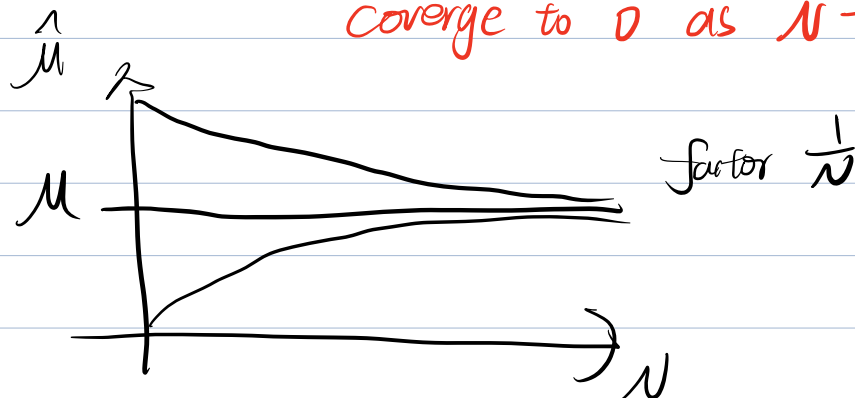$$= \frac{1}{N^2} E(\sum_i \sum_j (x_i - \mu)(x_j - \mu))$$

$i = j \Rightarrow E((x_i - \mu)^2) = \sigma^2$

$i \neq j \Rightarrow E[(x_i - \mu)(x_j - \mu)]$
$= E(x_i - \mu) E(x_j - \mu) = 0$

$$= \frac{1}{N^2}(N\sigma^2) = \frac{\sigma^2}{N} = Var(\hat{\mu})$$

converge to 0 as $N \to \infty$



factor $\frac{1}{N}$

Gaussian Variance (PS 2-12)

$$E(\hat{\sigma}^2) = \frac{N-1}{N} \sigma^2 \Rightarrow Bias(\hat{\sigma}^2) = -\frac{1}{N}\sigma^2 \neq 0$$

to make it unbiased:

$$\hat{\sigma}^2 = \frac{N}{N-1} \hat{\sigma}^2 = \frac{1}{N-1} \sum_{N} (x_i - \mu)^2$$
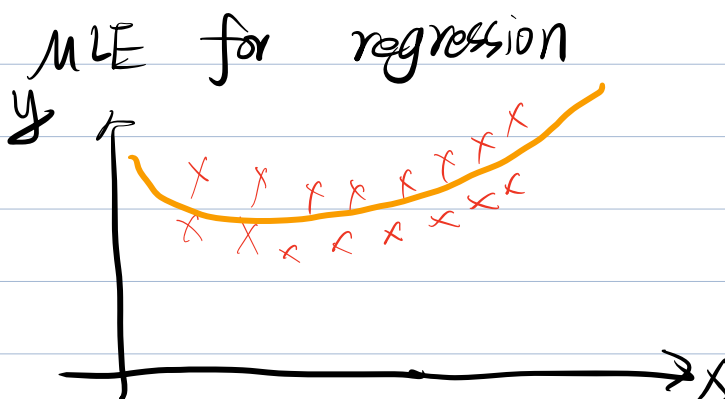
---

Important Asymptotic Properties of MLE

1) consistent: As $N \to \infty$, the estimate
converges to the true value
Asymptotically unbiased

2) efficient: achieves Cramer-Rao
Lower Bound (CRLB) as $N \to \infty$
- CRLB is a theoretical bound on the
variance of any unbiased estimator for
a given $p(x|\theta)$
- i.e. no unbiased estimator can
get lower variance.

---

MLE for regression



$x \in R$ input
$y \in R$ output
learn $f(x)$

Consider a $k^{th}$ order polynomial

$$\bar{f}(x;\theta) = \sum_{d=0}^{k} x^d \theta_d = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^k \end{bmatrix}^T \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{bmatrix} = \phi(x)^T \theta$$

Observe a noisy output:

$$y = f(x,\theta) + \varepsilon$$

noise $\varepsilon \sim N(0,\sigma^2)$ i.i.d.

random variable     deterministic     random variable

equivalently,

$$p(y|x,\theta) = N(y|f(x,\theta), \sigma^2)$$

Given dataset $\{(x_i, y_i)\}_{i=1}^{N}$, estimate $\theta$
using MLE

$$\hat{\theta} = \arg\max_{\theta} \sum_{i=1} \log p(y_i|x_i,\theta)$$

$$= \vdots$$

$$= \arg\min_{\theta} \sum_{i} (y_i - f(x_i,\theta))^2$$

$$= \arg\min_{\theta} \|y - \bar{\Phi}^T\theta\|^2, \quad \bar{\Phi} = [\phi(x_1),\dots,\phi(x_N)]$$

$$, y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Notes:
   ▷ MLE is more general than least square

2) assumptions are explicit
    i) Gaussian noise
    (ii) $\mu = 0$, $\sigma^2$ variance (fixed)
    (iii) noise is i.i.d.

3) MLE can describe other least square
    formulations:

change the
regression to
fit the tasks
    i) weighted LS (PS 2.8)
      (different Var)

↓
change the
noise distribution
    iii) regularised LS (lec. 3)
    iii) Lp-norm    (PS 2.9)

eg: non-negtve error ⇒ gamma
    non-neg integer error ⇒ poisson