

Refresh:

- In generative learning, the goal is to estimate

$$P(X, Y) = P(X|Y) P(Y).$$

- After that, we can use Bayes classifier that return $\underset{y}{\operatorname{argmax}} P(y|x)$

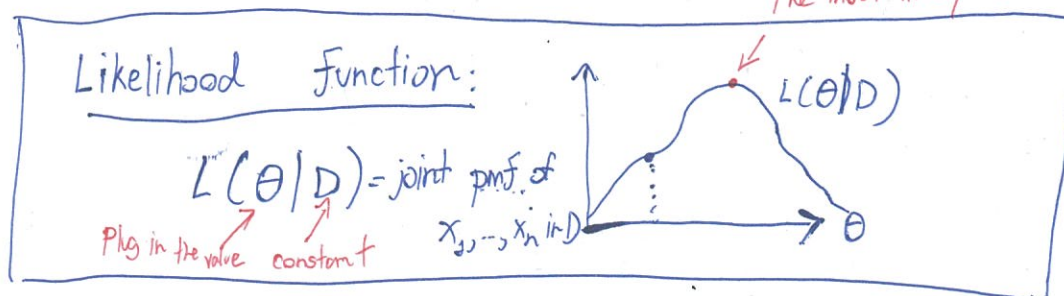
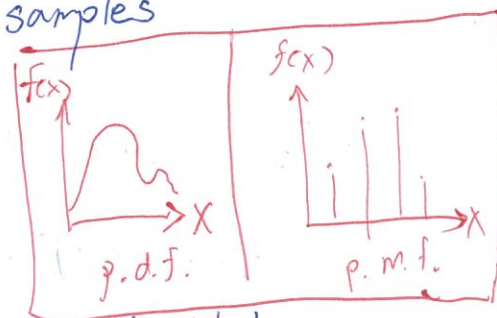
Estimating Probability from Data:

- Let θ be a parameter of the modelling distribution. $f(x|\theta)$
p.m.f. / p.d.f.

- The data set D contains n random samples from $f(x|\theta)$.

- Recall that the n samples are i.i.d.

- Likelihood estimation: Try to guess what is the most likely value that θ could be given the data D we have observed.



$$L(\theta|D = \{(x_1, y_1), \dots, (x_n, y_n)\}) = f(x_1, y_1 | \theta) \dots f(x_n, y_n | \theta)$$

$$= f(x_1, y_1 | \theta) \dots f(x_n, y_n | \theta) \quad (\text{Each sample is i.i.d.})$$

$$= \prod_{i=1}^n f(x_i, y_i | \theta)$$

Maximum Likelihood Estimation (MLE):

- Principle: find $\theta_{\max} = \arg \max_{\theta} (L(\theta|D))$
- The pmf. $f(D|\theta)$ ~~is the~~ represents the modelling distribution assumption

Simple scenario: freethrows:

- Let's say we attempt $n=6$ ^{i.i.d.} freethrows and $n_H=5$ make it in.

$$D = \{ \underset{x_1}{1}, 0, 1, 1, 1, \underset{x_n}{1} \}$$

Intuity, we should guess

$$P(X=1) = \frac{5}{6}$$

- Each $x_i \sim P(X)$ \leftarrow To estimate
r.v. denoting the ~~number of~~ successes at freethrows

- To estimate $P(X)$, we assume $x_i \sim f(D|\theta) = \text{Bin}(n, \theta)$

$$\Rightarrow f(D|\theta) = \binom{n}{n_H} \theta^{n_H} (1-\theta)^{n-n_H}$$

$$\Rightarrow L(\theta|D) = \binom{n}{n_H} \theta^{n_H} (1-\theta)^{n-n_H}$$

- In MLE, we will try to find θ_{\max}

$$\theta_{\max} = \arg \max_{\theta} \left(\binom{n}{n_H} \theta^{n_H} (1-\theta)^{n-n_H} \right)$$

$$= \arg \max_{\theta} \left(\log \left(\binom{n}{n_H} \theta^{n_H} (1-\theta)^{n-n_H} \right) \right)$$

$$= \arg \max_{\theta} \left(\underbrace{n_H \log \theta + (n-n_H) \log (1-\theta)}_{\log L(\theta|D)} \right)$$

- To find θ_{\max} , we ~~not~~ consider the first derivative

$$\frac{dL}{d\theta} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \frac{dL}{d\theta} = \frac{d(n_H \cdot \log \theta + (n - n_H) \cdot \log(1 - \theta))}{d\theta}$$
$$= \frac{n_H}{\theta} + \frac{n - n_H}{1 - \theta} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \frac{n_H}{\theta} = \frac{n - n_H}{1 - \theta}$$

$$\Rightarrow n_H(1 - \theta) = (n - n_H)(\theta)$$

$$\Rightarrow n_H - n_H \theta = \cancel{n - n_H} (n - n_H)(\theta)$$

$$\Rightarrow n \theta = n_H \Rightarrow \theta = \frac{n_H}{n} = \frac{5}{6}$$

Recap of MLE:

- Step 1: Make an explicit modeling assumption about what type of distribution the data was sampled from.

- Step 2: Set the parameter of the distribution so that the data observed is as likely as possible.

Simple Scenario II: Salmon or Mackinac

- For each type of fishes, we assume the ^{Gaussian} distribution

$$\Rightarrow x_i \sim P(x)$$

r.v. denotes the fish length

Our intuition is we compute the mean and s.d.
 \Rightarrow We have the gaussian distribution

- To estimate $P(X)$, we assume $x_i \sim \underbrace{f(x|\theta)}_{\text{p.d.f.}}$,
 where $\theta = \{\mu, \sigma\}$, and

$$f(x|\theta) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

- We define the likelihood function

$$\begin{aligned} L(\theta|D) &= f(D|\theta) \\ &= f(x_1|\theta) \cdot f(x_2|\theta) \cdots f(x_n|\theta) \\ &= \prod_{i=1}^n f(x_i|\theta) \end{aligned}$$

- Now, let's solve for θ_{\max} that maximizes $L(\theta|D)$

$$\theta_{\max} = \underset{\theta}{\operatorname{argmax}} \left(\prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2} \right)$$

$$= \underset{\theta}{\operatorname{argmax}} \left(\sum_{i=1}^n \log \left(\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2} \right) \right)$$

$$= \sum_{i=1}^n \left(-\frac{1}{2\sigma^2} (x_i - \mu)^2 \cdot \log \left(\frac{1}{\sigma\sqrt{2\pi}} e^{\frac{2}{2}} \right) \right)$$

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

$$\log \left(\prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \right) = \sum_{i=1}^n \log \left(\frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \right)$$

$$= \sum_{i=1}^n \log \left(\frac{1}{\sigma \sqrt{2\pi}} \right) + \log \left(e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \right)$$

$$= \sum_{i=1}^n \cancel{\log 1} - \log(\sigma \sqrt{2\pi}) + \log \left(e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \right)$$

$$= \sum_{i=1}^n -\log(\sigma) - \log(\sqrt{2\pi}) + \left(-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right)$$

$$= -n \log(\sigma) - n \log(\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$\log(L) = -\frac{n}{2} \log(\sigma^2) - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$\frac{\partial \log(L)}{\partial \sigma}$

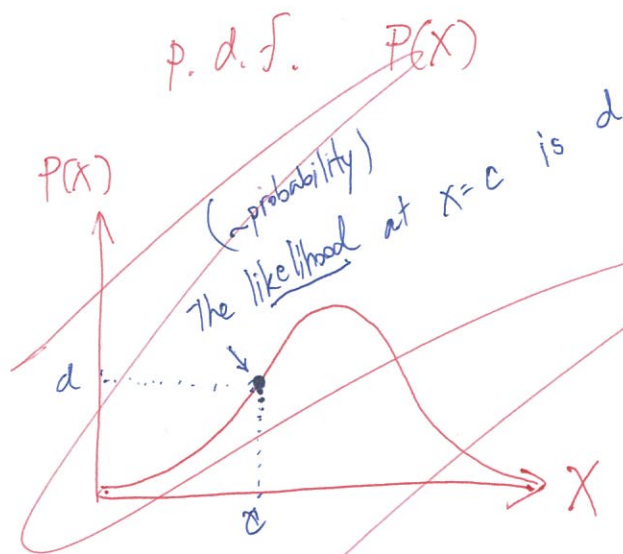
$$\frac{\partial \log(L)}{\partial \mu} = \frac{\partial}{\partial \mu} \left(\cancel{\frac{-x_i^2}{2\sigma^2} + \frac{2\mu x_i}{2\sigma^2} - \frac{\mu^2}{2\sigma^2}} \right) \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right)$$

$$= \frac{\partial}{\partial \mu} \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i^2 - 2\mu x_i + \mu^2) \right)$$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^n -2x_i + 2\mu = -\frac{1}{2\sigma^2} (-2) \sum_{i=1}^n (x_i - \mu)$$

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

$$\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \stackrel{\text{set}}{=} 0 \rightarrow \mu_{\max} = \frac{1}{n} \sum_{i=1}^n x_i$$



$$\sum_{i=1}^n -\frac{1}{2} \left(\frac{x_i - \mu}{6} \right)^2 \log \left(\frac{e}{6\sqrt{2\pi}} \right) = -\frac{1}{2} \cdot \log \left(\frac{e}{6\sqrt{2\pi}} \right) \cdot \sum_{i=1}^n \left(\frac{x_i - \mu}{6} \right)^2$$

$$= -\frac{1}{2} [1 + \log(\sigma) + \log(2\pi)] \cdot \sum_{i=1}^n \left(\frac{x_i - \mu}{6} \right)^2$$

$$= -\frac{1}{2} + \frac{\log 6}{2} - \frac{\log(2\pi)}{2}$$

$$= -\frac{1}{2 \cdot 6^2} \left[\sum_{i=1}^n (x_i - \mu)^2 \right] \left[\log e - \log 6\sqrt{2\pi} \right]$$

$$= \left(-\frac{1}{2 \cdot 6^2} \sum_{i=1}^n (x_i - \mu)^2 \right) - \frac{\log 6}{2 \cdot 6^2}$$

$$\frac{\partial \log L}{\partial \sigma} = \frac{\partial}{\partial \sigma} \left(-n \log(\sigma) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2 \right)$$

Fact:
 $\frac{\partial \ln(x)}{\partial x} = \frac{1}{x}$

$$= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \frac{n}{\sigma} = \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2$$

$$\Rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\sigma_{\max} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2} \quad !!!$$

Note:

$\mu_{\max}, \sigma_{\max} \Rightarrow \theta_{\max}$
 and they are what
 our intuition is
 about!!!

Summary of MLE:

- MLE gives the explanation of the data we observed.
- If n is large and your choice of distribution is correct, then MLE finds the "true" parameters
- MLE can overfit the data if n is small, it works well when n is large
- If you don't have the correct model, the MLE can be terribly wrong

