motivation problem of MLE

Coin flipping Bernoulli r.v. = $\{0 = T, 1 = H\}$

MLE: $\hat{\pi} = \frac{1}{N} \sum_i x_i$

Suppose we see: $D = \{1, 1, 1, 0, 0, 0, 0\} \Rightarrow \hat{\pi} = \frac{3}{7}$

what if $D' = \{1, 1, 1\}$ only $\hat{\pi} = 1$ ? (we never see tails)

This is an example of <u>overfitting</u>. (not enough samples to get a good estimate of the parameter)

what we can do?

- use our knowledge: we know $\pi \approx \frac{1}{2}$ for most coins and we incorperate this knowledge to our estimate of $\pi$.

Bayesian Param Estimation.

— treat $\theta$ as a <u>r.v.</u>

— Framework

—training set $D = \{x_1, \dots, x_N\}$

—prob density given parameter $\theta$: $p(x_i | \theta)$

— <u>prior</u> distribution on parameter $\theta$, $p(\theta)$ (added)
(encode prior beliefs about $\theta$, eg: $\pi \approx \frac{1}{2}$)

— posterior dist. of $\theta$ given data $D$.

$$p(\theta | D) = \frac{p(D|\theta) \cdot p(\theta)}{\int p(D|\theta) p(\theta) d\theta} \Rightarrow \text{density functions}$$

updated.

— predictive dist. — likelihood of new $x_*$ given data $D$.

how likely is the $\theta$ (weight) already fixed in the formula.

$$p(x_* | D) = \int p(x_* | \theta) p(\theta | D) d\theta$$

average over all $\theta$, weighted by posterior $p(\theta|D)$

"allow different explanations of data"

compared to $\underset{\text{pure determined by data}}{\underline{ML}}$, bayes is influenced by prior

$\Rightarrow$ (problem: how to get prior)

Example : Gaussian (known variance)

prior on $\mu$: $p(\mu) = N(\mu | \mu_0, a_0^2)$ ← prior are given

Likelihood of $x$: $p(x|\mu) = N(x | \mu, a^2)$

Dataset : $D = \{x_1, \cdots, x_N\}$

Calculate posterior

$$p(\mu|D) = \frac{[\prod_i^N p(x_i|\mu)]p(\mu)}{\int [\prod_i^N p(x_i|\mu)]p(\mu)d\mu}$$

← product of gaussian.

← doesn't depend on $\mu$.

※ Just look at numerator w.r.t. $\mu$, then normalize later.

product of Gaussian

can swap $(x-\mu)^2 = (\mu-x)^2$

$$N(x|a, A) \cdot N(x|b, B) = N(a | b, A+B) \, N(x|c, C)$$

$$C = \frac{1}{\frac{1}{A} + \frac{1}{B}} \Rightarrow \frac{1}{C} = \frac{1}{A} + \frac{1}{B} \qquad c = C(\frac{a}{A} + \frac{b}{B})$$

first 2 terms:

$$p(x_1|\mu) \cdot p(x_2|\mu) = N(\mu|x_1, a^2) \cdot N(\mu|x_2, a^2)$$

$$= N(x_1|x_2, 2a^2) \, N(\mu | \tilde{\mu}_2, \tilde{\sigma}_2^2)$$

$$\begin{cases} \dfrac{1}{\tilde{\sigma}_2^2} = \dfrac{1}{a^2} + \dfrac{1}{a^2} = \dfrac{2}{a^2} \\[2mm] \tilde{\mu}_2 = \dfrac{\sigma^2}{2}(\dfrac{x_1}{a^2} + \dfrac{x_2}{a^2}) = \dfrac{1}{2}(x_1 + x_2) \end{cases}$$

$$\Rightarrow p(x_1|\mu) \, p(x_2|\mu) \propto N(\mu | \tilde{\mu}_2, \tilde{\sigma}_2^2) \quad \text{(throw away the constant)}$$

first 3 terms

$$N(\mu | \tilde{\mu}_2, \tilde{\sigma}_2^2) \, N(x_3|\mu, a^2) \left( \propto N(\mu | \tilde{\mu}_3, \tilde{\sigma}_3^2) \right)$$

$$\text{precision} \leftarrow \boxed{\frac{1}{\tilde{\sigma_3}^2}} = \frac{1}{\tilde{\sigma}^2} + \frac{1}{\sigma^2} = \frac{2}{\sigma^2} + \frac{1}{\sigma^2} = \frac{3}{\sigma^2}$$

$$\tilde{\mu_3} = \frac{\sigma^2}{3} \left( \frac{\frac{1}{2}(x_1 + x_2)}{\frac{\sigma^2}{2}} + \frac{x_3}{\sigma^2} \right) = \frac{1}{3} (x_1 + x_2 + x_3)$$

$\Rightarrow$ reduction

first $N$ terms

$$\prod_{i}^{N} p(x_i | \mu) \propto N(\mu | \tilde{\mu}_n, \tilde{\sigma}_n)$$

$$\begin{cases} \tilde{\mu}_n = \frac{1}{N} \sum_{i}^{N} x_i = \boxed{\hat{\mu}_{ML}} \\ \tilde{\sigma}_\mu^2 = \frac{\sigma^2}{N} \end{cases}$$

turn out to be MLE res.

Add prior

mul constant

$$N(\mu | \tilde{\mu}_n, \tilde{\sigma}_n^2) \, N(\mu | \mu_0, \tilde{\sigma_0}^2) \propto N(\mu | \hat{\mu}_n, \hat{\sigma}_n^2)$$

$$\frac{1}{\hat{\sigma}_n^2} = \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \Rightarrow \hat{\sigma}_n^2 = \frac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}}$$

$$\hat{\mu}_n = \frac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}} \left( \frac{\hat{\mu}_{ML}}{\frac{\sigma^2}{N}} + \frac{\mu_0}{\sigma_0^2} \right) \quad \begin{array}{c} \text{mul} \\ \frac{\sigma_0^2 \sigma^2}{\sigma_0^2 \sigma^2} \end{array}$$

$$\hat{\mu}_n = \boxed{\frac{N \sigma_0^2}{\sigma^2 + N \sigma_0^2}} \hat{\mu}_{ML} + \boxed{\frac{\sigma^2}{\sigma^2 + N \sigma_0^2}} \cdot \mu_0$$

$$\alpha \qquad \qquad 1 - \alpha$$

Finally, $p(\mu | D) = N(\mu | \hat{\mu}_n, \hat{\sigma}_\mu^2)$

⭐ two constants $\Rightarrow$ cancel

# What does it mean ?

interpret between MLE $s_d$ and prior $\mu_0$

---

**Data Size**

**mean**
$$N = 0 \implies \alpha = 0 \implies \hat{\mu}_n = \mu_0$$
$$N \to \infty \implies \alpha = 1 \implies \hat{\mu}_n = \mu_{ML}$$

**variance**
$$N = 0 \implies \hat{\sigma}_n^2 = \sigma_0^2$$
$$N \to \infty \implies \hat{\sigma}_n^2 \to 0 \quad \leftarrow \text{Converge to single value.}$$

---

$$\sigma_0^2 \ll \sigma^2 \implies \alpha = 0 \implies \hat{\mu}_n = \mu_0$$

$\uparrow$

small prior $\longrightarrow$ ( strong belief compared to noise
$$\implies \text{use our belief} )$$

$$\sigma_0^2 \gg \sigma^2 \implies \alpha = 1 \implies \hat{\mu}_n = \hat{\mu}_{ML}$$

( weak belief $\to$ use MLE. )

$$\sigma^2 = \sigma_0^2 \implies \alpha = \frac{N}{N+1} \implies \hat{\mu}_n = \frac{1}{N+1}(N \cdot \hat{\mu}_{ML} \cdot \mu_0)$$
$$= \frac{1}{N+1}(\sum_v x_v + \mu_0)$$

add a virtual Sample
at $\mu_0$, then compute the mean

- for large $N$, the v.s. does not matter
- for small $N$, move the posterio towards $\mu_0$

[This is a form of regularization]

predictive distribution.

$$p(\mu \mid D) = N(\mu \mid \hat{\mu}_n, \hat{\sigma}_n^2)$$

$$p(x \mid \mu) = N(x \mid \mu, \alpha^2)$$

$$p(x \mid D) = \int p(x \mid \mu) \cdot p(\mu \mid D) d\mu$$

$$= \int N(\mu \mid x, \alpha^2) N(\mu \mid \hat{\mu}_n, \hat{\sigma}_n^2) d\mu$$

$$= \int N(x \mid \hat{\mu}_N, \alpha^2 + \hat{\sigma}_n^2) \underbrace{N(\mu \mid \cdots, \cdots) d\mu}_{\text{normalized to 1}}$$

<span style="color:red">integration over $\mu$</span>

$$p(x \mid D) = N(x \mid \hat{\mu}_N, \hat{\sigma}_n^2 + \alpha^2)$$

<span style="color:red">Same mean as posterior    variance of paramet~ M|D (uncertainty)    (uncertainty) due to noise observation</span>

Maximize a Posteriori (MAP)

Avoid calculating the denominator of Bayes' Rule

$$\int p(D \mid \theta) \, p(\theta) \, d\theta \quad \text{---- difficult.}$$

<span style="color:red">Solu. pick the $\theta$ with largest posterior prob.</span>

$$\hat{\theta}_{MAP} = \underset{\theta}{argmax} \; p(\theta \mid D)$$

$$= \underset{\theta}{argmax} \; \frac{p(D \mid \theta) \, p(\theta)}{\int p(D \mid \theta) \cdot p(\theta) d\theta} \Rightarrow \text{Constant w.r.t } \theta$$

$$= argmax \; p(D \mid \theta) \, p(\theta)$$

$$\hat{\Theta}_{MAP} = \underset{\theta}{argmax} \underbrace{\log p(D|\theta)}_{\text{data } \mathcal{L} \text{ for MLE}} + \underbrace{\log p(\theta)}_{\substack{\text{regularization} \\ \text{(belief)}}}$$

Example    Gaussian

$$\hat{\mu}_{MAP} = \underset{\mu}{argmax}\, p(\mu|D) = \underset{\mu}{argmax}\, N(\mu|\hat{\mu}_n, \hat{\sigma}_n^2)$$

$$\hat{\mu}_{MAP} = \hat{\mu}_n$$

Approximate posterior as a delta function:

$$P(\mu|D) \approx \delta(\mu - \hat{\mu}_n)$$
$$P(x|D) \approx p(x|\hat{\mu}_n) = N(x|\hat{\mu}_n, \sigma^2)$$

Bayesian Regression
  Same setup as before:
$$
\left.
\begin{array}{l}
\mathcal{BLR} \\
f(x) = \phi(x)^T \theta \\
y = f(x) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)
\end{array}
\right\} \Rightarrow
\begin{array}{l}
p(y|x,\theta) \\
= N(y|f(x), \sigma^2)
\end{array}
$$

Introduce prior on $\theta$:
$$p(\theta) = \mathcal{N}(\theta \mid 0, \alpha I)$$

↳ scaled identity
covariance matrix

MAP estimate

$$\hat{\theta} = \underset{\theta}{\arg\max} \; \log p(\theta \mid \theta) + \log p(\theta)$$

$$= \underset{\theta}{\arg\max} \; \sum_i \log p(y_i \mid x_i, \theta) + \log p(\theta)$$

tutorial $\vdots$

$$= \underset{\theta}{\arg\min} \; \| y - \bar{\Phi}^T \theta \|^2 + \lambda \| \theta \|^2$$

$$\hat{\theta} = (\bar{\Phi}\bar{\Phi}^T + \lambda I)^+ \bar{\Phi} y$$

↗ add some constant to eigenvalues

↖ constant ← controls regularization
$\lambda = 0 \Rightarrow L\text{-}S$
(depend on $a^2$ and $\alpha$)

ridge regression
- regularized least square
- shrinkage
- weight decay

regularize covariance matrix to prevent inverting an ill-conditioned matrix

(add $\lambda$ to all the eigenvalues of $\bar{\Phi}\bar{\Phi}^T$)

$$\begin{bmatrix} & \leftarrow \text{ridge}(+\lambda) \\ & \diagdown \\ & \end{bmatrix}$$

(sidenote:

$$(A + \alpha I)x = \alpha x$$
$$Ax = \lambda x$$
$$Ax + \alpha I x = \lambda x$$
$$Ax = (\lambda - \alpha I)x$$