# Lecture 3 - Bayesian Parameter Estimation

## Problem MLE

Coin 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 we see $D = \{1, 1, 1\}$ only $\Rightarrow \hat{\pi} = \frac{3}{3} = 1$ ?

This is <u>unreasonable</u>! We can only see H from this coin.
(we never see tails!)

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

• use our knowledge: we know $\pi \approx \frac{1}{2}$ for most coins.
Incorporate this knowledge into our estimate of $\pi$.

---

## Bayesian Param Estimation

- treat $\Theta$ as a <u>r.v.</u>

- <u>Framework</u>
  - training set $D = \{x_1, \dots, x_N\}$
  - prob. density given parameter $\Theta$: $p(x_i | \Theta)$
  - <u>prior</u> distribution on parameter $\Theta$: $p(\Theta)$
    (encodes prior beliefs about $\Theta$, eg. $\pi \approx \frac{1}{2}$)
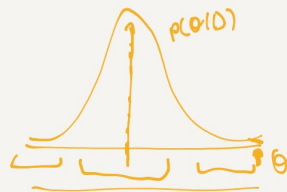
- <u>posterior dist.</u> of $\Theta$ given data D:
$$p(\Theta | D) = \frac{p(D|\Theta) p(\Theta)}{\int p(D|\Theta) p(\Theta) d\Theta} \quad \text{(Bayes' Rule)}$$

- <u>predictive dist.</u> - likelihood of new $x_*$ given data D,
$$p(x_* | D) = \int p(x_* | \Theta) p(\Theta | D) d\Theta$$

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

"allow different explanations of the data"

Example: Gaussian (known variance)

prior on $\mu$: $p(\mu) = N(\mu \mid \mu_0, \sigma_0^2)$ ← known

likelihood of $x$: $p(x \mid \mu) = N(x \mid \mu, \sigma^2)$ ← known

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

Calculate the posterior

$$p(\mu \mid D) = \frac{\left[\prod_i p(x_i \mid \mu)\right] p(\mu)}{\int \left[\prod_i p(x_i \mid \mu)\right] p(\mu) \, d\mu}$$

← product of Gaussians.

← doesn't depend on $\mu$. (constant) wrt $\mu$

☆ Just look at numerator wrt $\mu$, then normalize later.

first 2 terms:

$p(x_1 \mid \mu) p(x_2 \mid \mu) = N(x_1 \mid \mu, \sigma^2) N(x_2 \mid \mu, \sigma^2)$

$= N(\mu \mid x_1, \sigma^2) N(\mu \mid x_2, \sigma^2)$

$\underset{\text{"x"} \;\; a \;\; A}{\uparrow\uparrow\uparrow} \qquad \underset{\text{"x"} \;\; b \;\; B}{\uparrow\uparrow\uparrow}$

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

$C = \frac{1}{\frac{1}{A} + \frac{1}{B}} \quad \Rightarrow \frac{1}{C} = \frac{1}{A} + \frac{1}{B}$

$c = C\left(\frac{a}{A} + \frac{b}{B}\right)$

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

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

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

$p(x_1 \mid \mu) p(x_2 \mid \mu) \propto N(\mu \mid \tilde{\mu}_2, \tilde{\sigma}_2^2)$  (throw away the constant factor)

---

first 3 terms

$N(\mu \mid \tilde{\mu}_2, \tilde{\sigma}_2^2) N(x_3 \mid \mu, \sigma^2) \propto N(\mu \mid \tilde{\mu}_3, \tilde{\sigma}_3^2)$

$\underset{\text{precision}}{\underbrace{\dfrac{1}{\tilde{\sigma}_3^2}}} = \dfrac{1}{\tilde{\sigma}_2^2} + \dfrac{1}{\sigma^2} = \dfrac{2}{\sigma^2} + \dfrac{1}{\sigma^2} = \dfrac{3}{\sigma^2}$

$\tilde{\mu}_3 = \dfrac{\sigma^2}{3}\left(\dfrac{\frac{1}{2}(x_1+x_2)}{\frac{\sigma^2}{2}} + \dfrac{x_3}{\sigma^2}\right) = \dfrac{1}{3}(x_1 + x_2 + x_3)$

$\vdots$

first N terms:

$$\prod_i p(x_i \mid \mu) \propto N(\mu \mid \tilde{\mu}_n, \tilde{\sigma}_n^2) \quad \begin{cases} \tilde{\mu}_n = \dfrac{1}{N}\sum_i x_i = \hat{\mu}_{ML} \\[2mm] \tilde{\sigma}_n^2 = \dfrac{\sigma^2}{N} \end{cases}$$

× prior

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

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

$\hat{\mu}_n = \dfrac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}}\left(\dfrac{\hat{\mu}_{ML}}{\frac{\sigma^2}{N}} + \dfrac{\mu_0}{\sigma_0^2}\right) \quad \searrow \dfrac{\sigma_0^2 \sigma^2}{\sigma_0^2 \sigma^2}$

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

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

# What does it mean?

$$\hat{\mu}_n = \underbrace{\frac{N\sigma^2}{\sigma^2 + N\sigma_0^2}}_{\alpha} \hat{\mu}_{MC} + \underbrace{\frac{\sigma^2}{\sigma^2 + N\sigma_0^2}}_{1-\alpha} \mu_0$$

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

interpolating btwn MLE soln & prior $\mu_0$.

## Dataset size

$N = 0 \Rightarrow \alpha = 0 \Rightarrow \hat{\mu}_n = \mu_0$ ← no data, use prior.

$N \to \infty \Rightarrow \alpha = 1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{MC}$ ← lots of data, use MLE

## variance

$N = 0 \Rightarrow \hat{\sigma}_n^2 = \sigma_0^2$ ← prior uncertainty

$N \to \infty \Rightarrow \hat{\sigma}_n^2 \to 0$ ← converges to a single value.

---

$\sigma_0^2 \ll \sigma^2 \Rightarrow \alpha = 0 \Rightarrow \hat{\mu}_n = \mu_0$ (strong belief compared to noise → use our belief.)

$\sigma_0^2 \gg \sigma^2 \Rightarrow \alpha = 1 \Rightarrow \hat{\mu}_N = \hat{\mu}_{MC}$ (weak belief → use MLE)

---

$\sigma^2 = \sigma_0^2 \Rightarrow \alpha = \frac{N}{N+1} \Rightarrow \hat{\mu}_n = \frac{N}{N+1}\hat{\mu}_{MC} + \frac{1}{N+1}\mu_0$

$$= \frac{1}{N+1}\left(N\hat{\mu}_{MC} + \mu_0\right)$$

$$= \frac{1}{N+1}\left(\sum_i x_i + \mu_0\right)$$

"add a virtual sample at $\mu_0$, then "compute mean"

- for large $N$, the v.s. doesn't matter.
- for small $N$, moves the posterior towards $\mu_0$.

[ This is a form of regularization ]

# Predictive distribution

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

$$p(x | D) = \int p(x|\mu)p(\mu|D)\,d\mu$$

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

$$= \int \underbrace{N(x|\hat{\mu}_N, \sigma^2 + \hat{\sigma}_n^2)}_{\text{no } \mu} \underbrace{N(\mu|\cdots, \cdots)}\,d\mu$$

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

Same mean as posterior

variance of parameter $\mu | D$ (uncertainty)

uncertainty due to noisy obs.

$\sigma_0^2 \ll \sigma^2: \quad 1 - \alpha = \frac{\sigma^2}{\sigma^2 + N\sigma_0^2} = 1$

$\Rightarrow \alpha = 0$

# Maximum a Posteriori (MAP)

Avoid calculating the denominator of Bayes' Rule:

$$\int p(D|\theta)\, p(\theta)\, d\theta \quad \ldots \quad \text{difficult for many cases.}$$

Soln: pick the $\theta$ w/ largest posterior probability.

$$\hat{\theta}_{MAP} = \underset{\theta}{\arg\max}\ p(\theta|D)$$

$$\left.\begin{array}{c}\text{Bayes' Rule}\end{array}\right.$$

$$= \underset{\theta}{\arg\max}\ \frac{\overbrace{p(D|\theta)}^{\text{data like.}}\ \overbrace{p(\theta)}^{\text{prior}}}{\int p(D|\theta)\, p(\theta)\, d\theta}\quad \left.\begin{array}{l}\text{constant wrt } \theta\\ (\text{not a function of } \theta)\end{array}\right.$$

$$= \underset{\theta}{\arg\max}\ p(D|\theta)\, p(\theta)$$

$$\boxed{\hat{\theta}_{MAP} = \underset{\theta}{\arg\max}\ \underbrace{\log p(D|\theta)}_{\text{data LL for MLE}} + \underbrace{\log p(\theta)}_{\text{regularization}}}$$

## Example: Gaussian
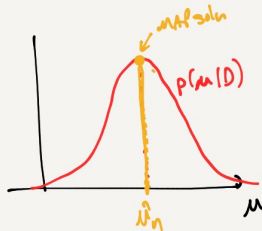
$$\hat{\mu}_{MAP} = \underset{\mu}{\arg\max}\ p(\mu|D) = \underset{\mu}{\arg\max}\ N(\mu|\hat{\mu}_n, \hat{\sigma}_n^2)$$

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

Approximate posterior as a delta function:

$$p(\mu|D) \approx \delta(\mu - \hat{\mu}_n)$$

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



MAP soln

$p(\mu|D)$

$\hat{\mu}_n$

$\mu$

# Bayesian Regression

Same setup as before:

$$x \in \mathbb{R}$$
$$f(x) = \phi(x)^T \theta \quad \leftarrow \theta \in \mathbb{R}^d$$
$$y = f(x) + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

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

Introduce prior on $\theta$: $p(\theta) = N(\theta\,|\,0, \alpha I)$

- identity matrix
- $\mathbb{R}^d$
- zero mean vector
- scaled identity covariance matrix

## MAP estimate

$$\hat{\theta} = \underset{\theta}{\arg\max}\ \log p(D|\theta) + \log p(\theta)$$

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

$$\vdots$$

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

constant

$$\boxed{\hat{\theta} = (\Phi\Phi^T + \lambda I)^{-1} \Phi y}$$

- ridge regression
  - regularized LS
  - Tikhonov regularization
  - shrinkage
  - weight decay

controls regularization

$\lambda = 0 \Rightarrow$ LS

$\lambda = \frac{\sigma^2}{\alpha}$ (see tutorial)

regularize the covariance matrix to prevent inverting an ill-conditioned matrix

(adds $\lambda$ to all the eigenvalues of $\Phi\Phi^T$)



ridge $(+\lambda)$