

Lecture 3: Bayesian Parameter Estimation

Problem w/ MLE

Model a coin as a Bernoulli r.v. $\{0=T, 1=H\}$

$$\text{MLE: } \hat{\pi} = \frac{1}{n} \sum_{i=1}^n x_i$$

Suppose we see: $D = \{1, 1, 1, 0, 0, 0, 1\} \Rightarrow \hat{\pi} = \frac{4}{7} \checkmark$

What if we see $D = \{1, 1, 1\} \Rightarrow \hat{\pi} = \frac{3}{3} = 1$

\Rightarrow This unreasonable $p(x=0) = 0 \Rightarrow$ ^{tails} never happens

\Rightarrow an example overfitting (too little data)

How to incorporate our knowledge about coins into our estimate of π ? e.g. $\pi = \frac{1}{2}$ typically.

Bayesian Parameter Estimation

- treat θ parameter as a r.v.

- Framework

- training set $D = \{x_1, \dots, x_n\}$

- = prob function: $p(x_i | \theta)$

- prior distribution on θ : $p(\theta) \Leftarrow$
(encode our beliefs about θ)

- posterior distribution of θ given D

$$p(\theta | D) = \frac{p(D | \theta) p(\theta)}{\int p(D | \theta) p(\theta) d\theta} \quad (\text{Bayes' Rule})$$

- predictive distribution - the likelihood of a new sample x_* given the data D .

$$p(x_* | D) = \int \underbrace{p(x_* | \theta)}_{\text{likelihood given a } \theta} \underbrace{p(\theta | D)}_{\text{posterior of } \theta} d\theta$$

"average over all θ weighted by its posterior probability"

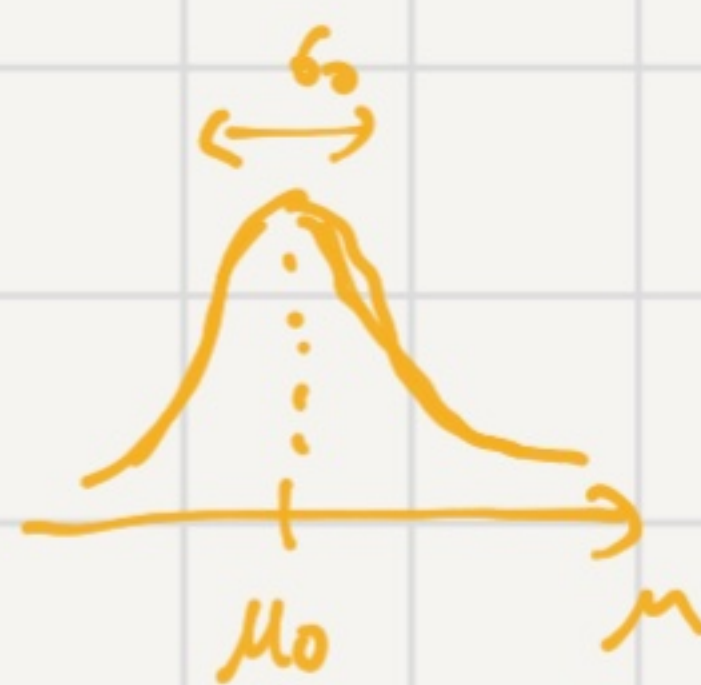
"allow different explanations of the data in terms of θ "

Example Gaussian (σ / known variance)

Dataset: $D = \{x_1, \dots, x_n\}$

likelihood: $p(x_i | \mu) = N(x_i | \mu, \sigma^2)$

prior: $p(\mu) = N(\mu | \mu_0, \sigma_0^2)$



Posterior

$$p(\mu | D) = \frac{p(D | \mu) p(\mu)}{\int p(D | \mu) p(\mu) d\mu} = \frac{\prod_{i=1}^n p(x_i | \mu) p(\mu)}{\int \prod_{i=1}^n p(x_i | \mu) p(\mu) d\mu}$$

not a function of $\mu \Rightarrow$ constant

$$= \frac{e^{f(\mu)}}{C} \leftarrow \text{what is } f(\mu)?$$

Just look at exponent of numerator & normalize later.

$$\begin{aligned} \log p(\mu | D) &\propto \sum_{i=1}^n \log p(x_i | \mu) + \log p(\mu) \\ &\propto \sum_{i=1}^n \left[-\frac{1}{2\sigma^2} (x_i - \mu)^2 \right] - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i^2 - 2\mu x_i + \mu^2) - \frac{1}{2\sigma_0^2} (\mu^2 - 2\mu_0\mu + \mu_0^2) \\ &\quad \underbrace{-2\mu \sum_{i=1}^n x_i + n\mu^2}_{n\hat{\mu}_{ML}} \\ &= -\frac{1}{2} \left(\underbrace{\frac{n}{\sigma^2} \mu^2}_{a} - 2 \underbrace{\left(\frac{n}{\sigma^2} \hat{\mu}_{ML} + \frac{1}{\sigma_0^2} \mu_0 \right) \mu}_{b} + \underbrace{\frac{1}{\sigma_0^2} \mu^2 + \frac{2}{\sigma_0^2} \mu_0 \mu}_{\text{constant}} \right) \\ &= -\frac{1}{2} \left[\underbrace{\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2}_a - 2 \underbrace{\left(\frac{n}{\sigma^2} \hat{\mu}_{ML} + \frac{1}{\sigma_0^2} \mu_0 \right) \mu}_b \right] \end{aligned}$$

Completing the square

$$ax^2 - 2bx + c = a(x - d)^2 + e$$

$$d = b/a$$

$$e = c - \frac{b^2}{a}$$

$$\begin{aligned} \Rightarrow \log p(\mu | D) &\propto -\frac{1}{2} \underbrace{\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)}_{\hat{\sigma}_n^{-2}} \left(\mu - \underbrace{\frac{1}{a} \left(\frac{n}{\sigma^2} \hat{\mu}_{ML} + \frac{1}{\sigma_0^2} \mu_0 \right)}_{\hat{\mu}_n} \right)^2 + \dots \\ &= -\frac{1}{2\hat{\sigma}_n^2} (\mu - \hat{\mu}_n)^2 + \text{const} \end{aligned}$$

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

$$\begin{aligned} \hat{\mu}_n &= \left(\frac{1}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right) \left(\frac{n}{\sigma^2} \hat{\mu}_{ML} + \frac{1}{\sigma_0^2} \mu_0 \right) \\ &= \left(\frac{\frac{n}{\sigma^2}}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right) \hat{\mu}_{ML} + \left(\frac{\frac{1}{\sigma_0^2}}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right) \mu_0 \end{aligned}$$

$$\hat{\mu}_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_{ML} + \left(\frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right) \mu_0$$

Thus

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

What does it mean?

$$\hat{\mu}_n = \underbrace{\left(\frac{n\sigma^2}{n\sigma^2 + \sigma_0^2} \right)}_{\alpha} \hat{\mu}_{ML} + \underbrace{\left(\frac{\sigma_0^2}{n\sigma^2 + \sigma_0^2} \right)}_{1-\alpha} \mu_0$$

$\hat{\mu}_{ML} \longleftrightarrow \mu_0$

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

Data size n:

$$\begin{aligned} n=0 &\Rightarrow \alpha=0 \Rightarrow \hat{\mu}_n = \mu_0 \\ n \rightarrow \infty &\Rightarrow \alpha=1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{ML} \end{aligned}$$

↗ smoothing b/w
↘ our prior belief μ_0
 & data-driven $\hat{\mu}_{ML}$

uncertainty

$$n=0 \Rightarrow \hat{\sigma}_n^2 = \sigma_0^2 \leftarrow \text{prior uncertainty}$$

$$n \rightarrow \infty \Rightarrow \hat{\sigma}_n^2 = 0 \leftarrow \text{converges to one value } \hat{\mu}_{ML}$$

Belief strength

$$\sigma_0^2 \ll \sigma^2 \Rightarrow \alpha=0 \Rightarrow \hat{\mu}_n = \mu_0 \Rightarrow \text{use prior}$$

↑
strong belief (small variance)

$$\sigma_0^2 \gg \sigma^2 \Rightarrow \alpha=1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{ML} \Rightarrow \text{use MLE}$$

↑
weak belief

$$\sigma_0^2 = \sigma^2 \Rightarrow \alpha = \frac{n\sigma^2}{n\sigma^2 + \sigma^2} = \frac{n}{n+1}$$

$$\Rightarrow \hat{\mu}_n = \frac{n}{n+1} \hat{\mu}_{ML} + \frac{1}{n+1} \mu_0$$

$$= \frac{1}{n+1} (n\hat{\mu}_{ML} + \mu_0)$$

$$= \frac{1}{n+1} \left(\sum_{i=1}^n x_i + \mu_0 \right)$$

$n+1$ terms in a sum

"add a virtual sample at μ_0 " then compute the mean.

- for large n , the v.s. doesn't matter.
- for small n , moves the MLE towards the prior mean.
- This is a form of regularization.

Predictive Distribution

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

likelihood: $p(x|\mu) = N(x|\mu, \sigma^2)$

predictive:

$$p(x|D) = \int p(x|\mu) p(\mu|D) d\mu = \int \underbrace{N(x|\mu, \sigma^2)}_{\dots e^{-(x-\mu)^2} = e^{-(\mu-x)^2}} N(\mu|\hat{\mu}_n, \hat{\sigma}_n^2) d\mu$$

$$= \int \underbrace{N(\mu|x, \sigma^2) N(\mu|\hat{\mu}_n, \hat{\sigma}_n^2)}_{\text{product of 2 Gaussians (PS 1-7)}} d\mu$$

product of 2 Gaussians (PS 1-7)

$$N(x|a, A) N(x|b, B) =$$

$$\underbrace{N(a|b, A+B)}_{\text{variance of posterior (uncertainty in } \mu)} \underbrace{N(x|c, C)}_{\text{variance due to noisy observations.}}$$

$$= \int \underbrace{N(x|\hat{\mu}_n, \sigma^2 + \hat{\sigma}_n^2)}_{\text{posterior mean}} \underbrace{N(\mu|\dots, \dots)}_{\text{variance due to noisy observations.}} d\mu$$

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

posterior mean

variance of posterior (uncertainty in μ)
variance due to noisy observations.

Maximum a Posteriori (MAP)

calculating $\int p(D|\theta) p(\theta) d\theta$ is difficult \Rightarrow approximations.

one solution: pick θ w/ highest posterior probability.

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(\theta|D)$$

$$= \underset{\theta}{\operatorname{argmax}} \frac{p(D|\theta) p(\theta)}{\int p(D|\theta) p(\theta) d\theta} \leftarrow \text{constant w.r.t } \theta$$

$$= \underset{\theta}{\operatorname{argmax}} p(D|\theta) p(\theta)$$

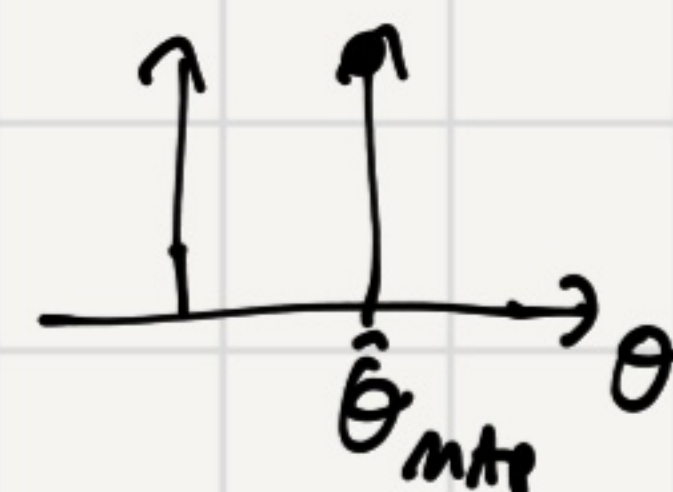
$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \log p(D|\theta) + \log p(\theta)$$

prior belief

posterior: $p(\theta|D) = \delta(\theta - \hat{\theta}_{\text{MAP}})$

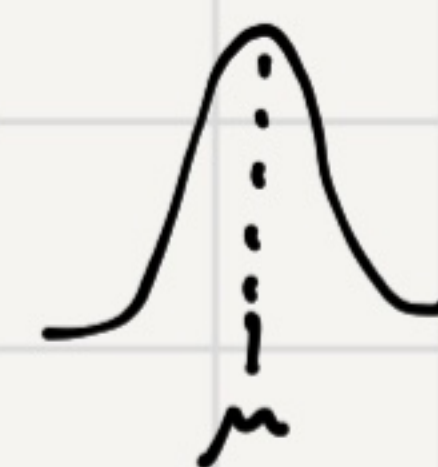
predictive: $p(x|D) = p(x|\hat{\theta}_{\text{MAP}})$

(ok if posterior is peaked, e.g. enough data)



Example: Gaussian

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



Bayesian Regression

similar setup:

$$x \in \mathbb{R}, \phi(x) \in \mathbb{R}^d, f(x, \theta) = \phi(x)^T \theta, \theta \in \mathbb{R}^d \text{ parameter vector}$$

$$y = f(x, \theta) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

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

New!

$$p(\theta) = N(\theta | 0, \alpha I)$$

parameter hyperparameter

Consider MAP estimate

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \underbrace{\log p(D|\theta)}_{\downarrow} + \underbrace{\log p(\theta)}_{\downarrow}$$

dropping constants

$$= \underset{\theta}{\operatorname{argmax}} -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \phi(x_i)^T \theta)^2 - \frac{1}{2\alpha} \|\theta\|^2 \cdot \sigma^2$$

$$= \underset{\theta}{\operatorname{argmin}} \underbrace{\sum_i (y_i - \phi(x_i)^T \theta)^2}_{\lambda} + \underbrace{\frac{\sigma^2}{\alpha} \|\theta\|^2}_{\lambda}$$

$$= \underset{\theta}{\operatorname{argmin}} \underbrace{\|y - \Phi^T \theta\|^2}_{\text{Squared error (same as L.S.)}} + \underbrace{\lambda \|\theta\|^2}_{\text{controls complexity of } \theta \text{ (prevent from taking extreme values)}}$$

hyperparameter controls tradeoff.

regularized L.S.
ridge regression
Tikhonov regularization
Shrinkage
weight decay

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

new λI term: regularizes the cov matrix to prevent inverting an ill-conditioned matrix.