

Lecture 3 Bayesian Parameter Estimation

problems w/ MLE

- Model a coin flip w/ Bernoulli r.v. $\{\theta = T, 1 = H\}$
- MLE: $\hat{\pi} = \frac{1}{N} \sum_{i=1}^N x_i$
- Suppose we see $D = \{1, 1, 1, 0, 0, 0, 1\} \Rightarrow \hat{\pi} = \frac{4}{7}$
- Suppose we see $D = \{1, 1, 1\} \Rightarrow \hat{\pi} = \frac{3}{3} = 1$
 $\pi = p(1) = 1$
 $1 - \pi = p(0) = 0 \leftarrow \text{probability } 0 \text{ of happening!}$
 (This is unreasonable \rightarrow we can never see tails!)
- This is an example of overfitting.
- We know that coins are usually fair ($\pi = \frac{1}{2}$).
 How to incorporate this knowledge into our estimator?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Bayesian Parameter Estimation

- treat parameters Θ as a r.v.
- Framework
 - training set: $D = \{x_1, \dots, x_N\}$
 - likelihood of data given parameter $p(x_i | \theta)$
 - prior distribution on parameters $p(\theta)$
 (encode our beliefs/knowledge about θ)
- posterior distribution of θ given the data D

$$p(\theta | D) = \frac{p(D|\theta) p(\theta)}{p(D)}, \quad p(D) = \int p(D|\theta) p(\theta) d\theta$$

(computing a distribution rather than a point estimate, like MLE)

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

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

(averaging over all θ weighted by the posterior)

"allow different explanations of the data"

Example: Gaussian (known variance)

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

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

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

posterior: $p(\mu | D) = \frac{p(D|\mu) p(\mu)}{p(D)}$

$$= \frac{\prod_{i=1}^N p(x_i|\mu) p(\mu)}{\int \prod_{i=1}^N p(x_i|\mu) p(\mu) d\mu} \quad \text{doesn't depend on } \mu$$

- look at the numerator first and see the form in terms of μ of the distribution, then normalize later.

$$\begin{aligned} \log p(\mu | D) &\propto \sum_{i=1}^N \log p(x_i|\mu) + \log p(\mu) \\ &= \sum_{i=1}^N \left[-\frac{1}{2\sigma^2} (x_i - \mu)^2 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 \right] \\ &\quad - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_0^2 \\ &\propto \sum_{i=1}^N \left[-\frac{1}{2\sigma^2} (x_i^2 - 2x_i\mu + \mu^2) - \frac{1}{2\sigma_0^2} (\mu^2 - 2\mu\mu_0 + \mu_0^2) \right] \\ &\propto -\frac{1}{2} \left[\frac{2}{\sigma^2} \sum_i x_i \mu + \frac{1}{\sigma^2} N \mu^2 + \frac{1}{\sigma_0^2} \mu^2 - \frac{2}{\sigma^2} \mu \mu_0 \right] \end{aligned}$$

$N \hat{\mu}_{ML}$

$$= -\frac{1}{2} \left[\left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{N}{\sigma^2} \hat{\mu}_{ML} + \frac{1}{\sigma_0^2} \mu_0 \right) \mu \right]$$

a b

$$a\mu^2 - 2b\mu$$

$$= a(\mu - \bar{\mu})^2 + e$$

completing the square: $\bar{\mu} = \frac{b}{a}$
 $e = -\frac{b^2}{a} = -a\bar{\mu}^2$

Assignment Project Exam Help
<https://powcoder.com>
 Add WeChat powcoder

$$= -\frac{1}{2} \frac{1}{\sigma_n^2} (\mu - \hat{\mu}_n)^2 + \text{const}$$

where $\sigma_n^2 = \frac{1}{\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}} = \frac{\sigma^2 \sigma_0^2}{N \sigma_0^2 + \sigma^2}$

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

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

Since the $\log p(\mu | D)$ is quadratic in μ

\Rightarrow Gaussian

$$p(\mu | D) = N(\mu | \hat{\mu}_n, \hat{\sigma}^2_n)$$
$$\text{let } \alpha = \frac{N\sigma_0^2}{\sigma^2 + N\sigma_0^2}$$
$$\hat{\mu}_n = \frac{N\sigma_0^2 \hat{\mu}_{ML} + \sigma^2 \mu_0}{N\sigma_0^2 + \sigma^2} = \alpha \hat{\mu}_{ML} + (1-\alpha)\mu_0$$
$$\hat{\sigma}^2_n = \frac{1}{\frac{1}{\sigma^2} + \frac{N}{\sigma_0^2}}$$

What's the interpretation?

Data size:

$$N=0 \Rightarrow \alpha=0 \Rightarrow \hat{\mu}_n = \mu_0 \quad (\text{no data, use prior})$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$N \rightarrow \infty \Rightarrow \alpha=1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{ML} \quad (\text{lots data, use MLE})$$

Smooth between MLE & prior.

$$N=0 \Rightarrow \hat{\sigma}^2_n = \sigma_0^2 \quad (\text{use prior uncertainty})$$

$$N \rightarrow \infty \Rightarrow \hat{\sigma}^2_n = 0 \quad (\text{converges to a single value, the MLE})$$

$$\hat{\sigma}^2_0 \ll \hat{\sigma}^2 \Rightarrow \alpha=0 \Rightarrow \hat{\mu}_n = \mu_0 \quad (\text{observations too noisy} \Rightarrow \text{use prior})$$

$$\hat{\sigma}^2_0 \gg \hat{\sigma}^2 \Rightarrow \alpha=1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{ML} \quad (\text{weak belief, use MLE})$$

$$\hat{\sigma}^2 = \hat{\sigma}_0^2 \Rightarrow \alpha = \frac{N}{N+1} \Rightarrow$$

$$\begin{aligned}\hat{\mu}_n &= \frac{N}{N+1} \hat{\mu}_{ML} + \frac{1}{N+1} \mu_0 \\ &= \frac{1}{N+1} (N \hat{\mu}_{ML} + \mu_0) \\ &= \frac{1}{N+1} \left(\sum_{i=1}^N x_i + \mu_0 \right)\end{aligned}$$

- adding 1 "virtual" sample at μ_0 , then compute the mean.
- for small N , the virtual sample has an effect.

• for large N , it has ^{little} effect.

This is a form of regularization.

Predictive distribution

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

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

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

N($\mu|x, \sigma^2$) ↓
Multiply 2 Gaussians (PSI-7)

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

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

↑ same as posterior ↑ observation noise uncertainty in μ .

Maximum a Posteriori (MAP)

With full Bayesian estimation we need

$$p(D) = \int p(D|\theta) p(\theta) d\theta \dots \text{usually hard to do.}$$

one solution: pick θ w/ highest posterior probability

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

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

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

↑ data log-likelihood (like MLE) ↑ prior (regularize the θ)

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Example: Gaussian

$$\hat{\mu}_{MAP} = \operatorname{argmax}_{\mu} p(\mu|D) = \hat{\mu}_n$$

predictive distribution: $p(x|\hat{\mu}_{MAP}) = N(x | \hat{\mu}_n, \sigma^2)$

Bayesian Regression

similar to before

$$\text{function: } f(x, \theta) = \Phi(x)^T \theta$$

$$\text{noisy obs: } y = f(x, \theta) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

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

$$\text{Data likelihood: } p(y|x, \theta) = N(y | f(x), \sigma^2)$$

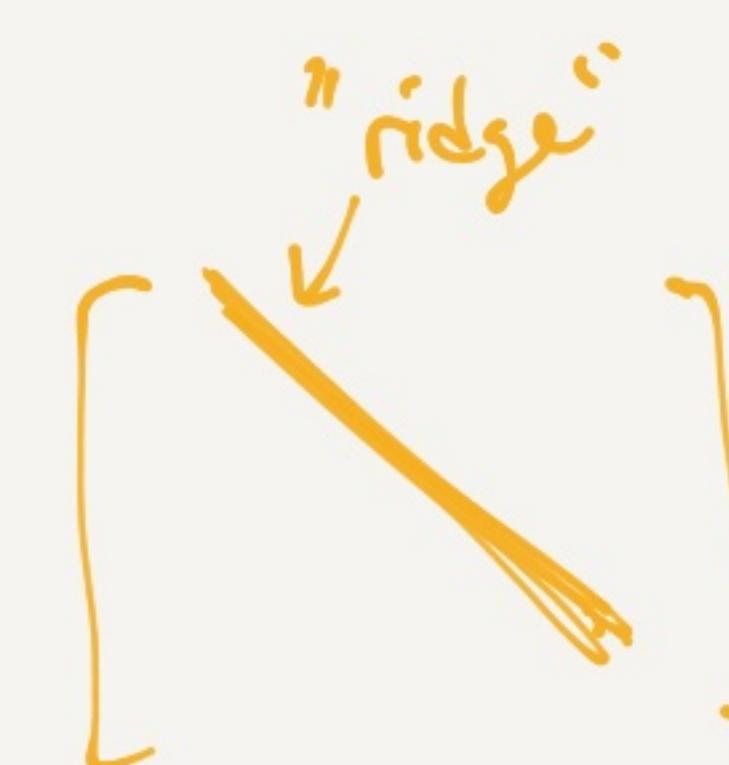
α = variance of the prior (known)

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

$$\text{L.S: } \hat{\theta} = (\Phi \Phi^T)^{-1} \Phi y$$

regularizes the covariance matrix $\Phi \Phi^T$

prevents ill-conditioned matrix & problems of inverting it.



MAP estimate θ

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

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log p(y_i|x_i, \theta) + \log p(\theta)$$

Same as tutorial

$$= \underset{\theta}{\operatorname{argmax}} -\frac{1}{2\sigma^2} \|y - \Phi^T \theta\|^2 + -\frac{1}{2\alpha} \|\theta\|^2 + \text{const.}$$

$$= \underset{\theta}{\operatorname{argmin}} \left(\frac{1}{2\sigma^2} \|y - \Phi^T \theta\|^2 + \frac{1}{2\alpha} \|\theta\|^2 \right) \cdot \sigma^2$$

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmin}} \|y - \Phi^T \theta\|^2 + \lambda \|\theta\|^2$$

hyperparameter

$\frac{\sigma^2}{2} = \lambda$
 controls the length of θ
 \Rightarrow controls the complexity of f .

$$f(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_p x^p$$

Assignment Project Exam Help
<https://powcoder.com>
 Add WeChat powcoder

Tutorial 3

PS 3.10

Bayesian Regression w/ Gaussian Prior.
parameter vector $\theta \in \mathbb{R}^P$

$$y = \Phi^T \theta + \epsilon \quad \leftarrow \text{noise } \epsilon \sim N(0, \Sigma)$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \Phi = \begin{bmatrix} 1 & \phi(x_1) & \dots & \phi(x_N) \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

Prior: $\theta \sim N(0, \Gamma)$

observation likelihood: $p(y|X, \theta) = N(y|\Phi^T \theta, \Sigma)$

a) Posterior of $\theta | X, y$

$$p(\theta | X, y) = \frac{p(y | X, \theta) p(\theta)}{p(y | X)}$$

(Bayes Rule)

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Look at the log of numerator to see the form of θ .

$$\log p(y | X, \theta) + \log p(\theta)$$

$$= -\frac{1}{2} \|y - \Phi^T \theta\|^2 \Sigma - \frac{1}{2} \log |\Sigma| - \frac{N}{2} \log 2\pi$$

$$- \frac{1}{2} \|\theta\|^2 \Gamma - \frac{1}{2} \log |\Gamma| - \frac{P}{2} \log 2\pi$$

$$\propto -\frac{1}{2} \left[y^T \Sigma^{-1} y - 2y^T \Sigma^{-1} \Phi^T \theta + \theta^T \Phi^T \Sigma^{-1} \Phi^T \theta \right]$$

$$+ \theta^T \Gamma^{-1} \theta$$

$$\propto -\frac{1}{2} \left[\theta^T (\Phi \Sigma^{-1} \Phi^T + \Gamma^{-1}) \theta - 2y^T \Sigma^{-1} \Phi^T \theta \right] - b^T$$

$$\theta^T A \theta - 2b^T \theta \Rightarrow (\theta - \hat{\mu})^T \hat{\Sigma}^{-1} (\theta - \hat{\mu})$$

P1.10 complete the square

$$f(x) = x^T A x - 2b^T x + c$$

$$= (x - d)^T A (x - d) + e$$

where $d = A^{-1} b$
 $e = c - b^T A^{-1} b$

$$\hat{\mu} = \underbrace{(\Phi \Sigma^{-1} \Phi^T + \Gamma^{-1})^{-1}}_{A^{-1}} \underbrace{\Phi \Sigma^{-1} y}_{b}$$

$$\hat{\Sigma} = \underbrace{(\Phi \Sigma^{-1} \Phi^T + \Gamma^{-1})^{-1}}_{A^{-1}}$$

$$\Rightarrow \log p(y | X, \theta) p(\theta) \propto -\frac{1}{2} (\theta - \hat{\mu})^T \hat{\Sigma}^{-1} (\theta - \hat{\mu})$$

since the log numerator is quadratic,

$$\Rightarrow p(\theta | X, y) = N(\theta | \hat{\mu}, \hat{\Sigma})$$

then it must be Gaussian
(constants will combine to form the Gaussian normalization)

$$\frac{1}{(2\pi)^{N/2} |\hat{\Sigma}|^{1/2}}$$

b) MAP estimate

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \ p(\theta | X, y) =$$

$$= \hat{\mu} = (\Phi \Sigma^{-1} \Phi^T + \Gamma^{-1}) \Phi \Sigma^{-1} y$$

$$\hat{\theta}_{MAP} = (\Phi \Sigma^{-1} \Phi^T + \Gamma^{-1})^{-1} \Phi \Sigma^{-1} y$$

↑
regularizer - helps make the inverse better conditioned.

$$\Sigma = \begin{bmatrix} \sigma^2 & & \\ & \ddots & 0 \\ 0 & \ddots & \sigma^2 \end{bmatrix} \Rightarrow \Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma^2} & & \\ & \ddots & 0 \\ 0 & \ddots & \frac{1}{\sigma^2} \end{bmatrix}$$

$$\Phi \Sigma^{-1} = \begin{bmatrix} \Phi_1 & \dots & \Phi_N \end{bmatrix} \begin{bmatrix} \frac{1}{\sigma^2} & & \\ & \ddots & 0 \\ 0 & \ddots & \frac{1}{\sigma^2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} \Phi_1 & \dots & \frac{1}{\sigma^2} \Phi_N \end{bmatrix}$$

large observation noise
 $\sigma_i^2 \rightarrow$ lower weight
 on this point (y_i, x_i)

⇒ Similar to weighted least squares.

⇒ Similar to regularized least squares

Bayesian regression combines them.

c) assume $\Gamma = \alpha I$, $\Sigma = \sigma^2 I$

$$\begin{aligned} \hat{\theta}_{MAP} &= (\Phi \Sigma^{-1} \Phi^T + \Gamma^{-1})^{-1} \Phi \Sigma^{-1} y \\ &= (\Phi \frac{1}{\sigma^2} \Phi^T + \frac{1}{\alpha} I)^{-1} \Phi \Sigma^{-1} y \\ &= (\Phi \Phi^T + \frac{\sigma^2}{\alpha} I)^{-1} \Phi y \\ \boxed{\hat{\theta}_{MAP} = (\Phi \Phi^T + \lambda I)^{-1} \Phi y} \end{aligned}$$

same as regularized least squares!

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \|y - \Phi^T \theta\|^2 + \lambda \|\theta\|^2 \\ &= \underset{\theta}{\operatorname{argmin}} y^T y - 2y^T \Phi^T \theta + \underline{\theta^T \Phi \Phi^T \theta} + \underline{\lambda \theta^T \theta} \end{aligned}$$

$$\frac{\partial}{\partial \theta} = -2\Phi^T y + \underbrace{2\Phi \Phi^T \theta + 2\lambda \theta}_{2(\Phi \Phi^T + \lambda I)\theta} = 0$$

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

Assignment Project Exam Help
<https://powcoder.com>
 Add WeChat powcoder

e) predictive distribution

Given an x_* , $f(x_*)$

$$f_* = f(x_*) = \theta^T \phi(x_*)$$

Note $\theta \sim p(\theta | x, y) = N(\theta | \hat{\mu}, \hat{\Sigma})$

$$p(f_* | x, y, x_*)$$

transformations of r.v. (P1.1)

$$y = Ax + b, x \text{ is r.v.}$$

$$E[y] = A E[x] + b$$

$$\text{cov}(y) = A \text{cov}(x) A^T$$

f_* is a r.v. that is a linear xfrm of $\theta|x_*$

$$\Rightarrow E[f_*] = \phi(x_*)^T \hat{\mu}$$

$$\Rightarrow \text{cov}(f_*) = \phi(x_*)^T \hat{\Sigma} \phi(x_*)$$

Since θ is Gaussian, then

$$f_* | x, y, x_* \sim N(f_* | \phi(x_*)^T \hat{\mu}, \phi(x_*)^T \hat{\Sigma} \phi(x_*))$$

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder