

Lecture 3 Bayesian Parameter Estimation.

Problem of MLE

Coin bernoulli R.V. $\{T=0, H=1\}$

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

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

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

This is unreasonable! we can only see head from this coin (we never see tails!)

This is an example of overfitting (not enough samples to get a good estimate of the parameter)

How to deal with this problem?

use our knowledge: we know $\pi \approx \frac{1}{2}$ for most coins

Incorporate this knowledge into our estimate of π

Framework:

- treat θ as a r.v.
- training set $D = \{x_1, \dots, x_N\}$
- probability distribution given parameter θ : $p(x_i | \theta)$
- prior distribution on parameter θ : $P(\theta)$
- encode prior beliefs about θ . e.g $\pi = \frac{1}{2}$

Posterior distribution of Θ given data D

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

Predictive distribution: likelihood of new X^* given data D.

$$P(X^*|D) = \int \underline{P(X^*|\theta) P(\theta|D)} d\theta$$

Average over all θ , weighted by Posterior $P(\theta|D)$

Allow different explanations of the data.

Example: Gaussian (Known Variance)

prior on μ : $P(\mu) = N(\mu | \mu_0, \sigma_0^2)$

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

We know μ_0 , σ_0^2 , and σ^2

Data set $D = \{x_1, \dots, x_N\}$

Calculate the posterior distribution

$$P(\mu|D) = \frac{\left[\prod_{i=1}^N P(D_i|\mu) \right] P(\mu)}{\int \left[\prod_{i=1}^N P(D_i|\mu) \right] P(\mu) d\mu} \quad \leftarrow \text{product of Gaussian}$$

doesn't depend on μ (consistent w.r.t. μ)

Just look at numerator w.r.t μ , then normalize later

First two terms:

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

$$\text{As } N(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

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

$$= N(\mu | x, \sigma^2) \text{ without meaning}$$

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

Product of Gaussian:

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

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

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

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

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

↓
constant

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

Throw away the constant factor

First Three items:

$$N(\mu | \tilde{\mu}_2, \tilde{6}_2^2) \cdot N(x_3 | \mu, 6^2) \propto N(\mu | \tilde{\mu}_3, \tilde{6}_3^2)$$
$$\left\{ \begin{array}{l} \frac{1}{\tilde{6}_3^2} = \frac{1}{\tilde{6}_2^2} + \frac{1}{6^2} = \frac{3}{6^2} \\ \tilde{\mu}_3 = \frac{1}{3} 6^2 \left(\frac{\tilde{\mu}_2}{\tilde{6}_2^2} + \frac{x_3}{6^2} \right) = \frac{6^2}{3} \left(\frac{2}{6^2} \cdot \frac{1}{2}(x_1 + x_2) + \frac{x_3}{6^2} \right) \\ = \frac{1}{3}(x_1 + x_2 + x_3) \\ \vdots \end{array} \right.$$

First N terms:

$$\frac{N}{1} P(x_i | \mu) \propto N(\mu | \tilde{\mu}_N, \tilde{6}_N^2)$$

$$\tilde{\mu}_N = \frac{1}{N} \sum x_i = \hat{\mu}_{ML} \quad \tilde{6}_N^2 = \frac{6^2}{N}$$

X prior distribution:

$$N(\mu | \tilde{\mu}_H, \tilde{6}_H^2) \propto N(\mu | \mu_0, 6_0^2) \propto N(\mu | \hat{\mu}_n, \hat{6}_n^2)$$

$$\frac{1}{\hat{6}_n^2} = \frac{1}{\tilde{6}_H^2} + \frac{1}{6_0^2} = \frac{1}{6^2} + \frac{1}{6_0^2} \Rightarrow \hat{6}_n^2 = \frac{6^2 6_0^2}{N 6_0^2 + 6^2}$$

$$\begin{aligned} \hat{\mu}_n &= \hat{6}_n^2 \left(\frac{\tilde{\mu}_H}{\tilde{6}_H^2} + \frac{\mu_0}{6_0^2} \right) = \frac{1}{\frac{6^2}{6^2} + \frac{1}{6_0^2}} \left(\frac{\tilde{\mu}_H}{6^2/N} + \frac{\mu_0}{6_0^2} \right) \\ &= \frac{N 6_0^2}{6^2 + N 6_0^2} \hat{\mu}_{ML} + \frac{6^2}{6^2 + N 6_0^2} \mu_0 \end{aligned}$$

Finally: $P(\mu | D) = N(\mu | \hat{\mu}_n, \hat{6}_n^2)$

What does it mean?

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

2

1-2

adjust between MLE solution AND prior μ_0

Dataset size:

$N=0 \Rightarrow \lambda=0 \Rightarrow \hat{\mu}_n = \mu_0$ NO DATA, USE PRIOR

$N \rightarrow \infty \Rightarrow \lambda=1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{ML}$ LOTS OF DATA USE MLE

Variance:

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

$N=0 \Rightarrow \hat{\sigma}_n^2 = \sigma^2$ PRIOR UNCERTAINTY

$N \rightarrow \infty \Rightarrow \hat{\sigma}_n^2 = 0$ CONVERGES TO A SINGLE VALUE

$\sigma^2 \ll \sigma_0^2 \Rightarrow \lambda=0 \Rightarrow \hat{\mu}_n = \mu_0$

Strong belief compared to noise \rightarrow USE OUR BELIEF

$\sigma^2 \ll \sigma_0^2 \Rightarrow \lambda=1 \Rightarrow \hat{\mu}_n = \hat{\mu}_{ML}$

weak belief \rightarrow USE MLE

$$\sigma^2 = \sigma_0^2 \Rightarrow \lambda = \frac{N}{N+1} \Rightarrow \hat{\mu}_n = \frac{1}{N+1} (N\hat{\mu}_{ML} + \mu_0)$$

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

add a virtual sample at μ_0 , then compute mean

Large N : the virtual sample doesn't matter

Small N : moves the posterior towards μ_0

THIS IS A FORM OF REGULARIZATION !

Predictive distribution

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

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

$$\begin{aligned} p(x|D) &= \int p(x|\mu) p(\mu|D) d\mu \\ &= \int N(x|\mu, \sigma^2) N(\mu|\hat{\mu}_n, \hat{\sigma}^2_n) d\mu \\ &= \int N(x|\hat{\mu}_n, \sigma^2 + \hat{\sigma}^2_n) N(\mu|..., ...) d\mu \end{aligned}$$

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

$\hat{\mu}_n$: Same mean as Posterior

$\hat{\sigma}^2_n$: Variance of Parameter $\mu|D$ (uncertainty)

σ^2 : Uncertainty due to noisy observe

Maximum a Posterior (MAP)

Avoid calculate the denominator of Baye's Rule

$$\int p(D|\theta) p(\theta) d\theta \text{ --- difficult for many cases}$$

Solution: pick the θ with largest Posterior Possibility

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} f(\theta | D) \rightarrow \text{data likelihood}$$

$$= \underset{\theta}{\operatorname{argmax}} \frac{p(D|\theta) p(\theta)}{\int p(D|\theta) p(\theta) d\theta} \rightarrow \text{Prior}$$

constant w.r.t θ - not a function of θ

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

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

data Likelihood for rule regularization

Example: Gaussian

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

$$= \underset{\mu}{\operatorname{argmax}} N(\mu | \hat{\mu}_n, \hat{\sigma}_n^2)$$

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

Approximate posterior as a delta function:

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

$$P(x | D) \approx P(x | \hat{G}_n) = N(x | \hat{\mu}_n, \hat{\Sigma})$$

Bayesian Regression:

Same setup as before

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

Introduce prior on θ : $P(\theta) = N(\theta | \theta_0, Q^{-1})$

$$\theta \in \mathbb{R}^d$$

$$\theta_0 \in \mathbb{R}^d$$

Q : scaled identity covariance matrix

MAP estimate

$$\hat{\theta} = \arg \max_{\theta} \log P(D|\theta) + \log P(\theta)$$

$$= \arg \max_{\theta} \sum_i p(x_i | \theta) + \log P(\theta)$$

• Ridge regression • regularized LS

• Tikhonov regularization • shrinkage

$$= \arg \min_{\theta} \|y - \Phi^T \theta\|^2 + \lambda \|\theta\|^2 \cdot \text{weight decay}$$

$\lambda = \frac{1}{2}$ controls regularization $\lambda = 0 \Rightarrow$ LS.

↓ regularize the covariance matrix

$\hat{\theta} = (\Phi \Phi^T + \lambda I)^{-1} \Phi y$ to prevent inverting an ill-conditioned matrix