

## Lecture 6: March 7

Lecturer: Anirban Bhattacharya &amp; Debdeep Pati

Scribes: Brian Kidd &amp; Yabo Niu

**Note:** *LaTeX template courtesy of UC Berkeley EECS dept & CMU's convex optimization course taught by Ryan Tibshirani.*

**Disclaimer:** *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

## 6.1 Three sufficient conditions

Let  $\mathcal{P}_n \subseteq \mathcal{P}$  denote the sieves and  $N(\epsilon_n, \mathcal{P}_n, d)$  denote the covering number for those sieves. Recall that there were three conditions for the theorem proving posterior contraction ( $\exists c_1, c_2 > 0$  constants):

$$\log N(\epsilon_n, \mathcal{P}_n, d) \leq c_1 n \epsilon_n^2 \quad (6.1)$$

$$\pi(p : D(p_0 || p) \leq \epsilon_n^2, V(p_0 || p) \leq \epsilon_n^2) \geq e^{-c_2 n \epsilon_n^2} \quad (6.2)$$

$$\pi(\mathcal{P}_n^C) \leq e^{-(c_2+4)n\epsilon_n^2} \quad (6.3)$$

If these three conditions hold, we can show that  $\mathbb{E}_0 \Pi[d(p, p_0) > M\epsilon_n \mid x_1, \dots, x_n] \rightarrow 0$  as  $n \rightarrow \infty$ .

## 6.2 Historical references

- Ghoshal, Ghosh, & van der Vaart: Convergence Rates of Posterior Distributions (2000, Annals of Statistics). This paper has some other useful variants, including the idea of using "shells" for quantifying prior probabilities on increasingly large circles, see Figure 6.1.

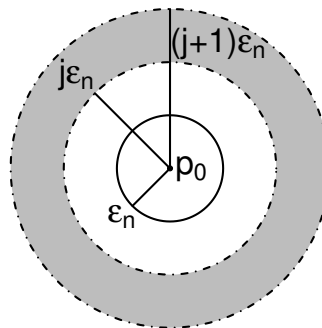


Figure 6.1: Breaking up the complement into shells.

- Shen & Wasserman: Rates of Convergence of Posterior Distributions (1999, Annals of Statistics)

- Ghoshal, Ghosh, & Ramamoorthi: Posterior consistency of Dirichlet mixtures in density estimation (1999, Annals of Statistics)
- Barron, Schervish, & Wasserman: The consistency of posterior distributions in nonparametric problems (1999, Annals of Statistics)
- Ghoshal & van der Vaart: Convergence rates of posterior distributions for noniid observations (2007, Annals of Statistics). This paper extends the ideas to non-iid (INID) cases such as regression. Some of the definition extensions include

$$p_0^{(n)} = \prod_{i=1}^n p_{0i}, \quad p = \prod_{i=1}^n p_i, \quad d(p_0^{(n)}, p^{(n)}) = \frac{1}{n} \sum_{i=1}^n d(p_{0i}, p_i).$$

### 6.3 Density Space & Parameter Space

In most practical situations,  $p \equiv p_\theta$ .

$$\begin{aligned} \Theta &\longrightarrow \mathcal{P} \\ \text{mapping: } \theta &\longmapsto p_\theta \end{aligned}$$

This map above doesn't have to be 1-1. Typically we place a prior distribution  $\Pi$  on  $\Theta$  which induces a prior  $\Pi$  on  $\mathcal{P}$ .

**Example 1:** DP mixture model  $x_1, \dots, x_n \mid p \stackrel{iid}{\sim} p$ .

$$p(x) = \sum_{h=1}^{\infty} \pi_h \cdot \frac{1}{\sigma_h} \phi\left(\frac{x - \mu_h}{\sigma_h}\right), \quad \text{location-scale mixtures of Gaussians}$$

$$\pi_h \geq 0, \quad \sum_{h=1}^{\infty} \pi_h = 1, \quad \text{a.s.}$$

$$\pi_h = \gamma_h \cdot \prod_{l < h} (1 - \gamma_l), \quad \text{stick-breaking representation}$$

$$\gamma_l \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad \mu_h \stackrel{iid}{\sim} \Pi_\mu, \quad \sigma_h \stackrel{iid}{\sim} \Pi_\sigma, \quad \text{priors}$$

$$\Pi = (\Pi_h), \quad \mu = (\mu_h), \quad \sigma = (\sigma_h), \quad p = p_\theta, \quad \theta = (\Pi, \mu, \sigma).$$

**Example 2:** (Regression)

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2).$$

The data is of the form  $(y_i, x_i)_{i=1}^n$ . Two common constructions for the function  $f$  are:

- (Possibly high-dimensional) linear regression:  $f(x_i) = x_i' \beta$ .
- Basis function regression (with bases  $\psi_j$ ):  $f(x_i) = \sum_{j=1}^K \theta_j \psi_j(x_i)$ .

### 6.4 Sieve construction

In most practical situations,  $p \equiv p_\theta$  (the densities are parameterized). We assume there is a measurable map from  $\Theta \rightarrow \mathcal{P}$ . Then a prior distribution  $\pi$  is placed on  $\Theta$ , which induces a prior  $\pi$  on  $\mathcal{P}$ . Denote the

distance metric on  $\Theta$  by  $d_\theta$ . For any subset  $\Theta_n \subseteq \Theta$ , define

$$\mathcal{P}_n = \{p_\theta : \theta \in \Theta_n\}$$

$$B(p_i, \epsilon_n) = \{p : d(p, p_i) < \epsilon_n\}.$$

We need a net on the sieves  $\mathcal{P}_n$  (i.e. we need  $p_1, \dots, p_N$  s.t.  $\mathcal{P}_n \subseteq \cup_{i=1}^N B(p_i, \epsilon_n)$ ). The idea is to find a  $\delta_n$ -net in the parameter space  $\Theta_n$  such that for any  $\theta, \theta' \in \Theta_n$ ,

$$d_\theta(\theta, \theta') \leq \delta_n \implies d(p_\theta, p_{\theta'}) \leq \epsilon_n$$

**Claim:** Define  $p_i \equiv p_{\theta_i}$ . Then  $\{p_1, \dots, p_N\}$  is an  $\epsilon_n$ -net of  $\mathcal{P}_n$ . For the complement probability,  $\pi(\mathcal{P}_n^C) \leq \pi(\Theta_n^C)$ .

## 6.5 Back to the density estimation example

$$x_1, \dots, x_n \mid p \stackrel{\text{iid}}{\sim} p,$$

$$p(x) = \sum_{h=1}^{\infty} \pi_h \frac{1}{\sigma_h} \phi\left(\frac{x_i - \mu_h}{\sigma_h}\right),$$

where  $\pi_h \geq 0$ ,  $\sum \pi_h = 1$  a.s. Note that  $\phi$  denotes the univariate normal pdf here, though the mixture model generalizes to multivariate normals and other distributions. For a prior on the  $(\pi_h)$ , we use the common stick-breaking prior:

$$\pi_h = \nu_h \prod_{\ell < h} (1 - \nu_\ell), \quad \nu_\ell \stackrel{\text{iid}}{\sim} \text{Beta}(1, \alpha).$$

To generalize to a Pitman-Yor prior, the Beta distribution should be defined with sequences  $\alpha_\ell, \beta_\ell$  for the parameters. For the means and standard deviations, any iid prior can be used (though later we'll see the need for the prior to have exponential tails in this example). Let  $\Pi_\mu, \Pi_\sigma$  denote these priors. Define the vectors  $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\sigma}$  to be the sequences  $\{(\pi_h), (\mu_h), (\sigma_h)\}$ , respectively, and

$$p \equiv p_\theta, \quad \theta = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\sigma}).$$

**Useful Inequalities:** Let  $q$  denote a density and  $\pi$  denote the weights for the mixture.

$$\left\| \sum_{h=1}^{\infty} \pi_h q_h - \sum_{h=1}^{\infty} \pi'_h q_h \right\|_{TV} \leq \sum_{h=1}^{\infty} |\pi_h - \pi'_h|$$

$$\left\| \sum_{h=1}^{\infty} \pi_h q_{1h} - \sum_{h=1}^{\infty} \pi_h q_{2h} \right\|_{TV} \leq \sum_{h=1}^{\infty} \pi_h \|q_{1h} - q_{2h}\|_{TV}$$

### 6.5.1 Sieve construction

The construction follows Dr. Pati's paper. This construction is useful, as it generalizes much better than other distributions.

$$\Theta_n = \left\{ (\Pi, \boldsymbol{\mu}, \boldsymbol{\sigma}) : \underbrace{\sum_{h=1}^{K_n} \Pi_h}_{\text{first } K_n \text{ weights make up most of the total mass.}} > (1 - \delta_n), \mu_h \in [-m_n, m_n], \sigma_h \in [u_n, v_n], h = 1, \dots, K_n \right\}.$$

The first  $K_n$  weights make up most of the mass, and both the means and variances are bounded.  $(m_n) \rightarrow \infty, (u_n) \rightarrow 0, (v_n) \rightarrow \infty$ . Let the distance metric be the total variation metric for the density space. Let  $p_\theta \equiv p_{\Pi, \mu, \sigma}$ . For  $\theta, \theta' \in \Theta_n$ , by triangle inequality,

$$\begin{aligned} \|p_\theta - p_{\theta'}\|_{\text{TV}} &\leq \|p_{\Pi, \mu, \sigma} - p_{\Pi', \mu, \sigma}\|_{\text{TV}} + \|p_{\Pi', \mu, \sigma} - p_{\Pi', \mu', \sigma'}\|_{\text{TV}} \\ &\leq \sum_{h=1}^{\infty} |\Pi_h - \Pi'_h| + \sum_{h=1}^{\infty} \Pi'_h \cdot \|N(\mu_h, \sigma_h^2) - N(\mu'_h, \sigma_h'^2)\|_{\text{TV}} \\ &\leq \sum_{h=1}^{K_n} |\Pi_h - \Pi'_h| + 2\delta_n + \sum_{h=1}^{K_n} \Pi'_h \cdot \frac{|\mu_h - \mu'_h|}{\min\{\sigma_h, \sigma'_h\}} + \delta_n \\ &\leq \sum_{h=1}^{K_n} |\Pi_h - \Pi'_h| + \frac{1}{u_n} \cdot \max_{h=1:K_n} \{|\mu_h - \mu'_h|\} + 3\delta_n. \end{aligned}$$

If we can make

$$\sum_{h=1}^{K_n} |\Pi_h - \Pi'_h| < \frac{\epsilon_n}{3}, \quad \frac{1}{u_n} \cdot \max_{h=1:K_n} |\mu_h - \mu'_h| < \frac{\epsilon_n}{3}, \quad \delta_n = \frac{\epsilon_n}{9},$$

then  $\|p_\theta - p_{\theta'}\|_{\text{TV}} < \epsilon_n$ .

**REMARKS:** If a space  $\Theta$  can be written as  $\Theta_a \otimes \Theta_b$  and  $\theta = (\theta_a, \theta_b)$ . Assume  $d(\theta, \theta') \leq d_a(\theta_a, \theta'_a) + d_b(\theta_b, \theta'_b)$ . Let  $\{\theta_{1a}, \dots, \theta_{Na}\}$  be an  $\epsilon_n/2$ -net of  $\Theta_a$  and  $\{\theta_{1b}, \dots, \theta_{Mb}\}$  be an  $\epsilon_n/2$ -net of  $\Theta_b$ . Then  $\{\theta_i = (\theta_{ia}, \theta_{jb}) : i = 1, \dots, N; j = 1, \dots, M\}$  is an  $\epsilon$ -net of  $\Theta$  with covering number  $\leq MN$ .

Let's look at the complement probability. By using Bonferroni inequality,

$$\begin{aligned} \Pi(\Theta_n^c) &= \Pi\left(\left\{\sum_{h=1}^{K_n} \Pi_h > (1 - \delta_n)\right\}^c \cup \bigcup_{h=1}^{K_n} \left\{\mu_h \in [-m_n, m_n]\right\}^c \cup \bigcup_{h=1}^{K_n} \left\{\sigma_h \in [u_n, v_n]\right\}^c\right) \\ &\leq \Pi\left(\left\{\sum_{h=1}^{K_n} \Pi_h > (1 - \delta_n)\right\}^c\right) + \Pi\left(\bigcup_{h=1}^{K_n} \left\{\mu_h \in [-m_n, m_n]\right\}^c\right) + \Pi\left(\bigcup_{h=1}^{K_n} \left\{\sigma_h \in [u_n, v_n]\right\}^c\right) \\ &\leq \Pi\left(\sum_{h=K_n+1}^{\infty} \Pi_h > \delta_n\right) + K_n \cdot \left[\Pi(\mu_1 \notin [-m_n, m_n]) + \Pi(\sigma_1 \notin [u_n, v_n])\right]. \end{aligned}$$

**Result/Verify:**

$$\begin{aligned} \sum_{h=K+1}^{\infty} \Pi_h &= 1 - \sum_{h=1}^K \Pi_h \\ &= 1 - \sum_{h=1}^K \left\{ \gamma_h \prod_{l=1}^{h-1} (1 - \gamma_l) \right\} \\ &= (1 - \gamma_1) - \gamma_2(1 - \gamma_1) - \gamma_3(1 - \gamma_1)(1 - \gamma_2) - \dots - \gamma_K(1 - \gamma_1) \cdots (1 - \gamma_{K-1}) \\ &= (1 - \gamma_1)(1 - \gamma_2) - \gamma_3(1 - \gamma_1)(1 - \gamma_2) - \dots - \gamma_K(1 - \gamma_1) \cdots (1 - \gamma_{K-1}) \\ &\dots \dots \dots \\ &= (1 - \gamma_1) \cdots (1 - \gamma_K) \end{aligned}$$

$$\text{Thus, } \sum_{h=K+1}^{\infty} \Pi_h = (1 - \gamma_1) \cdots (1 - \gamma_K).$$

By using the result above,

$$\begin{aligned}\Pi\left(\sum_{h=K+1}^{\infty} \Pi_h > \delta\right) &= \Pi\left(\sum_{h=1}^K \underbrace{-\log(1 - \gamma_h)}_{\text{Gamma}(1, \alpha)} < \log\left(\frac{1}{\delta}\right)\right) \\ &= \Pi\left(\text{Gamma}(K, \alpha) < \log\left(\frac{1}{\delta}\right)\right) \\ &= \text{exponentially small.}\end{aligned}$$