

Example Sheet 1 (of 3)

TBB/Mich 2017

Comments and corrections to t.berrett@statslab.cam.ac.uk. Starred* questions will be marked for the examples class.

[Notation: For a square-integrable function $g : \mathbb{R} \rightarrow \mathbb{R}$, define $R(g) = \int_{-\infty}^{\infty} g(x)^2 dx$; for a kernel K , define $\mu_2(K) = \int_{-\infty}^{\infty} x^2 K(x) dx$.]

1. Let $U_1, \dots, U_n \stackrel{iid}{\sim} U(0, 1)$, and let $Y_1, \dots, Y_{n+1} \stackrel{iid}{\sim} \text{Exp}(1)$. Writing $S_j = \sum_{i=1}^j Y_i$ for $j = 1, \dots, n+1$, show that

$$U_{(j)} \stackrel{d}{=} \frac{S_j}{S_{n+1}} \sim \text{Beta}(j, n-j+1),$$

for $j = 1, \dots, n$.

2.* (Hoeffding's inequality) (a) Let Y be a random variable with mean zero and $a \leq Y \leq b$. Use convexity to show that for every $t \in \mathbb{R}$, we have

$$\log \mathbb{E}(e^{tY}) \leq -\alpha u + \log(\beta + \alpha e^u),$$

where $u = t(b-a)$ and $\alpha = 1 - \beta = -a/(b-a)$. Using a second-order Taylor expansion about the origin, deduce that $\log \mathbb{E}(e^{tY}) \leq t^2(b-a)^2/8$.

(b) Now let Y_1, \dots, Y_n be independent with $\mathbb{E}(Y_i) = 0$ and $a_i \leq Y_i \leq b_i$ for $i = 1, \dots, n$. Use Markov's inequality to show that, for every $\epsilon > 0$, we have

$$\mathbb{P}\left(\left|\sum_{i=1}^n Y_i\right| > \epsilon\right) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

3. Let X_1, \dots, X_n be independent with distribution P on a measurable space $(\mathcal{X}, \mathcal{A})$, and let \hat{P}_n be the empirical measure of X_1, \dots, X_n ; thus $\hat{P}_n(A) = n^{-1} \sum_{i=1}^n \mathbb{1}_{\{X_i \in A\}}$ for $A \in \mathcal{A}$. Show that, for all $\epsilon > 0$ and $A \in \mathcal{A}$, we have

$$\mathbb{P}(|\hat{P}_n(A) - P(A)| > \epsilon) \leq 2e^{-2n\epsilon^2}.$$

4. (a) Let $X_1, \dots, X_n \stackrel{iid}{\sim} F$, and let \hat{F}_n denote their empirical distribution function. For $t_1 < \dots < t_k$, write down the distribution of

$$n(\hat{F}_n(t_1), \hat{F}_n(t_2) - \hat{F}_n(t_1), \dots, \hat{F}_n(t_k) - \hat{F}_n(t_{k-1}), 1 - \hat{F}_n(t_k)).$$

(b) Find the asymptotic distribution of $n^{1/2}(\hat{F}_n(t_1) - F(t_1), \dots, \hat{F}_n(t_k) - F(t_k))$.

5. (Continuation) We say a continuous process $(B_t)_{t \in [0,1]}$ is a *standard Brownian motion* on $[0, 1]$ if $B_0 = 0$, and if, for $0 \leq s_1 \leq t_1 \leq \dots \leq s_k \leq t_k \leq 1$, we have $(B_{t_1} - B_{s_1}, \dots, B_{t_k} - B_{s_k}) \sim N_k(0, \Sigma)$, where $\Sigma := \text{diag}(t_1 - s_1, \dots, t_k - s_k)$. The process $(W_t)_{t \in [0,1]}$ defined by $W_t = B_t - tB_1$ is called a *Brownian bridge*, or *tied-down Brownian motion*, because $W_0 = W_1 = 0$. Compute the distribution of $(W_{t_1}, \dots, W_{t_k})$.

[These last two questions suggest that “ $n^{1/2}(\hat{F}_n(t) - F(t)) \xrightarrow{d} W_{F(t)}$ as $n \rightarrow \infty$ ”. Care is required to make this statement and its proof precise.]

6. (a) Verify the algebraic identity

$$\phi_\sigma(x - \mu)\phi_{\sigma'}(x - \mu') = \phi_{\sigma\sigma'/(\sigma^2 + \sigma'^2)^{1/2}}(x - \mu^*)\phi_{(\sigma^2 + \sigma'^2)^{1/2}}(\mu - \mu'),$$

where $\mu^* = (\sigma'^2\mu + \sigma^2\mu')/(\sigma^2 + \sigma'^2)$, and $\phi_\sigma(x)$ is the $N(0, \sigma^2)$ density.

(b) Let X_1, \dots, X_n be independent $N(0, \sigma^2)$ random variables. Taking K to be the $N(0, 1)$ density, show that the mean integrated squared error of the kernel density estimate \hat{f}_h with kernel K and bandwidth h can be expressed exactly as

$$\text{MISE}(\hat{f}_h) = \frac{1}{2\pi^{1/2}} \left\{ \frac{1}{nh} + \left(1 - \frac{1}{n}\right) \frac{1}{(h^2 + \sigma^2)^{1/2}} - \frac{2^{3/2}}{(h^2 + 2\sigma^2)^{1/2}} + \frac{1}{\sigma} \right\}.$$

7. (Continuation) Now suppose that $h = h_n$ satisfies $h \rightarrow 0$ as $n \rightarrow \infty$ and $nh \rightarrow \infty$ as $n \rightarrow \infty$. Derive an appropriate asymptotic expansion of the *MISE* computed above, and deduce that the asymptotically optimal bandwidth with respect to the *MISE* criterion is given by

$$h_{AMISE} = \left(\frac{4}{3n}\right)^{1/5} \sigma.$$

Check that the same expression is obtained from the general formula for the asymptotically optimal bandwidth for a second-order kernel.

8. Let $X_1, \dots, X_n \stackrel{iid}{\sim} f$, where f'' is bounded. Write \tilde{f}_b for the histogram estimator of f with binwidth b . Assume $b = b_n \rightarrow 0$ and $nb \rightarrow \infty$ as $n \rightarrow \infty$. For $x \in \mathbb{R}$, let $I_b(x)$ denote the bin containing x and $p_b(x) = \mathbb{P}\{X_1 \in I_b(x)\}$ denote the bin probability. Show that

$$p_b(x) = bf(x) + \frac{1}{2}f'(x)[b^2 - 2b\{x - t_b(x)\}] + O(b^3)$$

as $n \rightarrow \infty$, where $t_b(x)$ is the left-hand endpoint of $I_b(x)$. Deduce that

$$\text{MSE}\{\tilde{f}_b(x)\} = \frac{f(x)}{nb} + \frac{1}{4}b^2 f'(x)^2 + f'(x)^2 \{x - t_b(x)\}^2 - bf'(x)^2 \{x - t_b(x)\} + O\left(\frac{1}{n} + b^3\right).$$

9. (Continuation) Assuming in addition that $R(f') < \infty$, argue informally that

$$\text{MISE}(\tilde{f}_b) = \frac{1}{nb} + \frac{1}{12}b^2 R(f') + o\left(\frac{1}{nb} + b^2\right).$$

Hence derive the AMISE optimal binwidth b_{AMISE} and find $\text{AMISE}(\tilde{f}_{b_{\text{AMISE}}})$.

10. (Scheffé's theorem) Let (f_n) be a sequence of densities and f be another density such that $f_n \rightarrow f$ almost everywhere. By integrating $g_n = f - f_n$ separately over $\{x : g_n(x) > 0\}$ and $\{x : g_n(x) \leq 0\}$ and using dominated convergence, show that

$$\int_{-\infty}^{\infty} |f_n(x) - f(x)| dx \rightarrow 0.$$

11.* Assume the standard conditions on f , h and K from lectures, and also that f'' is continuous with $R(f'') < \infty$. Use Fubini's theorem to show that $h \int_{-\infty}^{\infty} (K_h^2 * f)(x) dx = R(K)$.

Use the dominated convergence theorem to show that $(K_h * f)(x) \rightarrow f(x)$ for each $x \in \mathbb{R}$, and show that $\sup_{n \in \mathbb{N}} \sup_{x \in \mathbb{R}} (K_h * f)(x) < \infty$. Apply Scheffé's theorem to deduce that $\int_{-\infty}^{\infty} (K_h * f)^2(x) dx \rightarrow \int_{-\infty}^{\infty} f(x)^2 dx$.

Finally, deduce that

$$\int_{-\infty}^{\infty} \text{Var}\{\hat{f}_h(x)\} dx = \frac{1}{nh} R(K) + O(n^{-1}).$$

12. (Continuation) Show that $\int_{-\infty}^{\infty} [\mathbb{E}\{\hat{f}_h(x)\} - f(x)]^2 dx = h^4 \int_{-\infty}^{\infty} A_n^2(x) dx$, where

$$A_n(x) = \int_{-\infty}^{\infty} \int_0^1 (1-t) f''(x - thz) z^2 K(z) dt dz.$$

Apply Cauchy-Schwarz twice, firstly to the innermost integral with $(1-t)^{1/2}|z|K^{1/2}(z)$ as one term of the product, and secondly to the middle integral, and then use Fubini's theorem to evaluate the x -integral first, to show that

$$\int_{-\infty}^{\infty} A_n^2(x) dx \leq \frac{1}{4} R(f'') \mu_2^2(K)$$

for all n . Use dominated convergence to show that $A_n(x) \rightarrow \frac{1}{2} f''(x) \mu_2(K)$ for each $x \in \mathbb{R}$. Apply Fatou's lemma and combine the previous results to conclude that

$$\text{MISE}(\hat{f}_h) = \frac{1}{nh} R(K) + \frac{1}{4} h^4 R(f'') \mu_2^2(K) + o\left(\frac{1}{nh} + h^4\right).$$