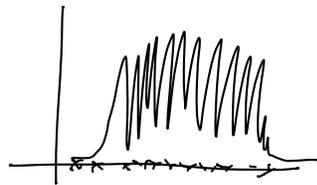
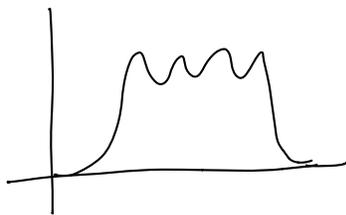


Density Estimation

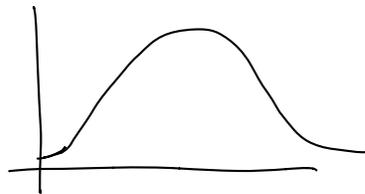
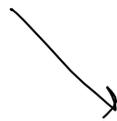
Setup: Observe data $X_1, \dots, X_n \sim F$, and the density is $\underline{f = F'}$.

Goal: Estimate f using as few assumptions as possible.

\Rightarrow Still a smoothing problem:



undersmoothed estimate



oversmoothed

One possible measure of the error is the L^2 error:

$$\begin{aligned}
 \text{Loss} = L &= \int (\hat{f}(x) - f(x))^2 dx \\
 &= \|\hat{f} - f\|^2 \\
 &= \underbrace{\|\hat{f}\|^2 - 2(\hat{f}, f)}_J + \|f\|^2 = J + C.
 \end{aligned}$$

$$\begin{aligned}
 (\hat{f}, f) &= \text{inner product of } \hat{f} \text{ with } f \\
 &= \int \hat{f}(x) f(x) dx \\
 &= \int \hat{f}(x) dF(x) \\
 &= E(\hat{f}(x))
 \end{aligned}$$

Goal is to estimate J .

As before, denote by $\hat{f}_{(-i)}$ the estimator obtained by leaving out x_i :

Def: CV estimate of the risk:

$$\underline{\underline{\hat{J} = \|\hat{f}\|^2 - \frac{2}{n} \sum \hat{f}_{(-i)}(x_i)}}$$

Histograms

Assume we are estimating f on $[0, 1]$, set $h = \frac{1}{m}$, then we have bins $B_1 = [0, h)$, $B_2 = [h, 2h)$, ..., $B_j = [(j-1)h, jh)$.

Denote by $Y_j = \# X_i$'s in bin j .

$\hat{p}_j = Y_j/n$. \leftarrow ^{estimate} probability of ending up in bin j

$p_j = \int_{B_j} f(x) dx$ \leftarrow true probability of landing in bin j .
 $= P(X \in B_j)$.

Histogram estimator: $\hat{f}(x) = \sum_{j=1}^m \frac{\hat{p}_j}{h} \mathbb{1}(x \in B_j)$.

\leftarrow Maybe a surprising factor?

Why not just \hat{p}_j ? \hat{f}

$|B_j| \cdot \frac{\hat{p}_j}{h} = \hat{p}_j \approx \int f(x)$

$$E(\hat{f}(x)) = \frac{E(\hat{p}_j)}{h} = \frac{p_j}{h} = \frac{1}{h} \int_{B_j} f(x) dx \approx \frac{1}{h} f(x) \cdot h = f(x).$$

Thm $E(\hat{f}(x)) = \frac{p_j}{h}$ for $x \in B_j$

$$\text{Var}(\hat{f}(x)) = \frac{p_j(1-p_j)}{n h^2}$$

and the risk can be computed as:

Then Assume that f' is "absolutely continuous" and

$\int (f')^2 < \infty$, then

$$R(\hat{f}, f) = \frac{h^2}{12} \int (f'(x))^2 dx + \frac{1}{nh} + O(h^2) + O\left(\frac{1}{n}\right)$$

and for fixed n , the minimum occurs at

$$h_* = \frac{1}{n^{1/3}} \left(\frac{6}{\int f'^2 dx} \right)^{1/3} \sim \frac{1}{n^{1/3}} \quad \text{and}$$

then $R(\hat{f}, f) \sim C \frac{1}{n^{2/3}}$.

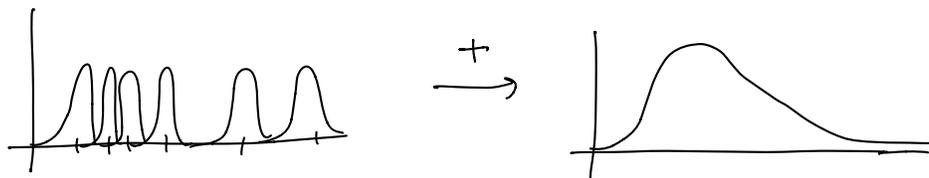
Remember: Since f is not known, minimize \hat{J} instead since it can actually be computed.

Kernel Density Estimator

If you had only one data point, x_i , what would you do?



Idea: Place a local kernel at each data point, and sum:



How wide should the kernel be?

Def: Kernel density estimator:

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x-x_i}{h}\right)$$

$$\int K = 1$$

$$\int xK = 0$$

$$\sigma_k^2 = \int x^2 K < \infty$$

Thm If f is continuous at x , and $h \rightarrow 0$, $nh \rightarrow \infty$,
 then $\hat{f}(x) \xrightarrow{P} f(x)$:

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\hat{f}(x) - f(x)| > \epsilon) = 0$$

Thm (6.28) Let $R(x) = \mathbb{E}(|f(x) - \hat{f}(x)|^2)$ be the risk
 at x . Then

$$R(x) = \frac{1}{4} \sigma_w^4 h^4 f''(x)^2 + \frac{f(x)}{nh} \int K^2(x) dx + o\left(\frac{1}{n}\right) + o(h^6).$$

Proof: $\mathbb{E}(\hat{f}(x)) = \mathbb{E}\left(\frac{1}{n} \sum_i \frac{1}{h} K\left(\frac{x-x_i}{h}\right)\right)$

$$= \mathbb{E}\left(\frac{1}{n} K\left(\frac{x-x_i}{h}\right)\right)$$

$$= \frac{1}{h} \int K\left(\frac{x-t}{h}\right) f(t) dt$$

$$= \int K(u) f(x-hu) du \quad , \quad \text{expand } f \text{ around } hu = 0.$$

$$= \int K(u) \left(f(x) - hu f'(x) + \frac{h^2 u^2}{2} f''(x) + \dots \right) du$$

$$= f(x) + \frac{1}{2} h^2 f''(x) \underbrace{\int u^2 K(u) du}_{\sigma_K^2} + \dots$$

assuming K is even.

$$\Rightarrow \text{bias} = \mathbb{E}(\hat{f}(x) - f(x)) = \frac{1}{2} h^2 \sigma_K^2 f''(x) + o(h^4)$$

$$\text{Similarly, } \text{var}(\hat{f}(x)) = \frac{f(x)}{nh} \int K^2(u) du + o\left(\frac{1}{n}\right).$$

$$\Rightarrow R = \text{bias}^2 + \text{variance}$$

Then for the optimal bandwidth,

$$\text{soln. } \frac{dR}{dh} = 0 \Rightarrow h \sim \frac{\sigma_k^2}{n^{1/5}}$$

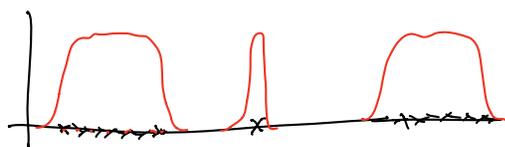
$$\Rightarrow R \sim O\left(\frac{1}{n^{4/5}}\right)$$

vs for the histogram $\therefore R \sim O\left(\frac{1}{n^{2/3}}\right)$

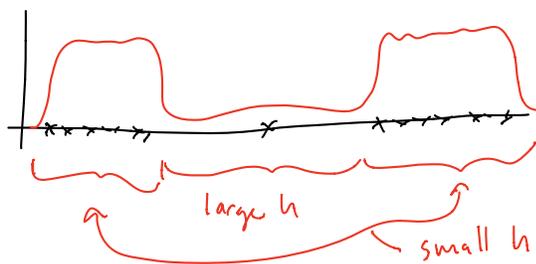
Note Assuming only that $\int (f'')^2 < \infty$, the rate of $\frac{1}{n^{4/5}}$ is the best that can be obtained (see Thm 6.31 in AONPS.)

Adaptive Methods Choose h locally depending on clustering of data and other considerations.

Ex:



uniformly small h



Multivariate Version

Same mathematical idea as in the one-dimensional case.

One option is to use what is known as a product kernel =

$$K_h^d(x_1, \dots, x_d) = \prod_{j=1}^d \frac{1}{h_j} K\left(\frac{x_j - x_j'}{h_j}\right)$$

K_h^d must satisfy kernel properties.

$$\begin{aligned} \text{then } \hat{f}(\vec{x}) &= \frac{1}{n} \sum_{i=1}^n K_h^d(\vec{x} - \vec{x}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \prod_{j=1}^d \frac{1}{h_j} K\left(\frac{x_j - x_{ij}}{h_j}\right) \end{aligned}$$

Risk can be estimated in the same way using multivariable Taylor series for f .

Curse of Dimensionality

If we want the risk $R \sim 0.1$ at $\vec{x} = 0$ for $f \sim \text{Normal}(0,1) \in \mathbb{R}^d$, using the optimal bandwidth then $n \sim c^d$!

d	n
1	4
2	19
4	223
8	43,700
9	187,000
10	842,000

Bootstrap (Ch. 8 in AoS) (C&B 10.1.4, 10.6.5) (AoS NPS Ch. 3).

Goal: Estimate standard errors and confidence sets for statistics.

Outline: Given $X_1, \dots, X_n \sim F$, statistic $T = g(X_1, \dots, X_n)$, want $\text{Var}_F(T)$.
 \curvearrowright dependence on unknown distribution F .

Ex: $T = \bar{X}$

$$\text{Var}_F = \frac{\sigma^2}{n}$$

$$\text{if } \text{var } X_i = \sigma^2$$

$$= \int (x - \mu)^2 dF(x)$$

function of F .

The idea of the bootstrap:

① Estimate $\text{Var}_F(T)$ with $\text{Var}_{\hat{F}}(T)$.

↖ \hat{F} put $1/n$ mass at every x_i .

② Use simulation to approximate $\text{Var}_{\hat{F}}(T)$.

→ In this example, step 2 is not needed because

$$\text{Var}_{\hat{F}}(T) = \frac{\hat{\sigma}^2}{n} \quad \text{when} \quad \hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2.$$

What is simulation? Drawing samples from some distribution, and computing averages.

Ex: Draw Y_1, \dots, Y_m from a distribution G , by the law of large numbers

$$\bar{Y} = \frac{1}{m} \sum_{j=1}^m Y_j \xrightarrow{\mathbb{P}} \mathbb{E}(Y) = \int y dG(y) \quad \text{as } m \rightarrow \infty.$$

Choosing m large enough means that $\bar{Y} \approx \mathbb{E}(Y)$, use this as an estimate for $\mathbb{E}(Y)$.

Also if h is some function with $\int h(y) dy < \infty$,

$$\text{then} \quad \frac{1}{m} \sum h(Y_j) \xrightarrow{\mathbb{P}} \mathbb{E}(h(Y)) = \int h(y) dG(y).$$

$$\underline{\text{Ex:}} \quad \frac{1}{m} \sum (Y_i - \bar{Y})^2 = \frac{1}{m} \sum Y_i^2 - \left(\frac{1}{m} \sum Y_i \right)^2$$

$$\xrightarrow{\mathbb{P}} \int y^2 dG(y) - \left(\int y dG(y) \right)^2$$

$$= \text{Var}(Y).$$

Bootstrap Variance Estimate

If we have data X_i , but F is unknown, then estimate F with \hat{F} , and draw from \hat{F} .

\Rightarrow Draw X_1^*, \dots, X_n^* from X_1, \dots, X_n with replacement.

\Rightarrow Compute $T^* = g(X_1^*, \dots, X_n^*)$

Note Some of the X_i^* will be duplicates.

Variance Algorithm

DO $i = 1, \dots, m$

Draw X_1^*, \dots, X_n^* from \hat{F}

Compute $T_i^* = g(X_1^*, \dots, X_n^*)$

COMPUTE
$$v_{\text{boot}} = \frac{1}{m} \sum_{j=1}^m (T_j^* - \bar{T}^*)^2$$

\uparrow
bootstrap estimate
of the variance

$\Rightarrow \hat{se} = \sqrt{v_{\text{boot}}}$

We can use exactly the same algorithm to estimate the variance of median, mode, or any other integrable statistic. $\int g < \infty$.

Bootstrap Confidence Intervals

Method 1 If T is approximately normal, e.g. an MLE, the T^* is also approximately normal (and so is v_{boot})

$$\text{Iterate: } \theta^{j+1} = \theta^j - \frac{l'(\theta^j)}{l''(\theta^j)}$$

Can show that if θ^j is "close enough" to root,
then $|\theta^{j+1} - \hat{\theta}| \sim |\theta^j - \hat{\theta}|^2$

quadratic convergence.

$ \theta^j - \hat{\theta} $	$ \theta^{j+1} - \hat{\theta} $
10^{-1}	10^{-2}
10^{-2}	10^{-4}
10^{-4}	10^{-8}
10^{-8}	10^{-16}

Notes

(1) θ^0 can usually be estimated by the method of moments estimator.

(2) l'' must be computed or approximated.

In the multivariate case, $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$

$H = \text{Hessian}$

and

$$l'(\vec{\theta}) = \nabla l(\vec{\theta}) \approx \begin{pmatrix} \frac{\partial l}{\partial \theta_1}(\theta_1^0) \\ \vdots \\ \frac{\partial l}{\partial \theta_k}(\theta_k^0) \end{pmatrix} + \begin{pmatrix} \frac{\partial^2 l}{\partial \theta_1^2} & \frac{\partial^2 l}{\partial \theta_1 \partial \theta_2} & \dots \\ \vdots & \ddots & \ddots \\ \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} \theta_1 - \theta_1^0 \\ \theta_2 - \theta_2^0 \\ \vdots \\ \theta_k - \theta_k^0 \end{pmatrix}$$

At $\hat{\theta}$, $\nabla l(\hat{\theta}) = \vec{0}$, so

$$\vec{\theta}^{j+1} = \vec{\theta}^j - H^{-1}(\vec{\theta}^j) (\nabla l(\vec{\theta}^j))$$

Multivariate
Newton's Method.

Stochastic Processes

or $X(t) = X_t$
 X_t is a r.v.

A stochastic process $\{X_t : t \in T\}$ is a collection of random variables indexed by t

- X_t takes values in the state space \mathcal{X}
- T is the index set (i.e. $\mathbb{R}, \mathbb{N}, \dots$)
- Ex include stock prices, weather, IID ~~square~~ X_1, \dots, X_n, \dots
- Recall: for X_1, \dots, X_n the joint density is given by

$$f(x_1, \dots, x_n) = f(x_1) f(x_2 | x_1) f(x_3 | x_1, x_2) \dots f(x_n | x_1, \dots, x_{n-1})$$
$$= \prod_{i=1}^n f(x_i | \text{past } i\text{'s})$$

Markov Chains

Def $\{X_n : n \in T\}$ is a Markov Chain

if $P(X_n = x | X_0, \dots, X_{n-1}) = P(X_n = x | X_{n-1})$

for all $n \in T$ and $x \in \mathcal{X}$.

$$\Rightarrow f(x_n | x_{n-1}, \dots, x_0) = f(x_n | x_{n-1})$$

$$\Rightarrow f(x_1, x_2, \dots, x_n) = f(x_1) f(x_2 | x_1) f(x_3 | x_2) \dots f(x_n | x_{n-1})$$

Questions to answer:

① When does a MC achieve "equilibrium"? Does it at all?

② Estimate parameters controlling the MC

③ Can we construct a MC that converges to a specified equilibrium? i.e. $X_n \rightsquigarrow F$, some given distribution

Transition Probability

Def: $p_{ij} = P(X_{n+1} = j \mid X_n = i)$ are the transition probabilities.

If p_{ij} does not depend on n , called a homogeneous MC.

The matrix P with elements $P_{ij} = p_{ij}$ is known as the transition matrix.

Two properties of the transition probabilities:

① $p_{ij} \geq 0$

② $\sum_j p_{ij} = 1$ (Typo in book).

↖ Each row of P is a prob. mass function.