# Extending the Basic Bootstrap (II)

Ex: Ten subjects were randomly assigned to each of two treatments. Find

(a) A bootstrap SE for $\bar{X}_1 - \bar{X}_2$.

(b) A bootstrap 95% CI for $\mu_1 - \mu_2$

(c) A bootstrap 95% CI for $\sigma_1^2 / \sigma_2^2$.

How would we do this?

Note: There are multiple possible approaches that one might use. A key thing for success of the bootstrap is to <u>mimic the original sampling as much as possible.</u>

# Two-Sample Example

| Treatment | Values |
|-----------|--------|
| 1 | 9  12  12  14  17 |
|   | 19  21  22  26  31 |
| 2 | 8  9  10  11  13 |
|   | 13  19  21  22  24 |

# Two methods based on Kernels

① Kernel density estimation — Allows us to

estimate the [probability density function] for a  pdf

distribution w/o 桶 assuming a parametric form

for the density.

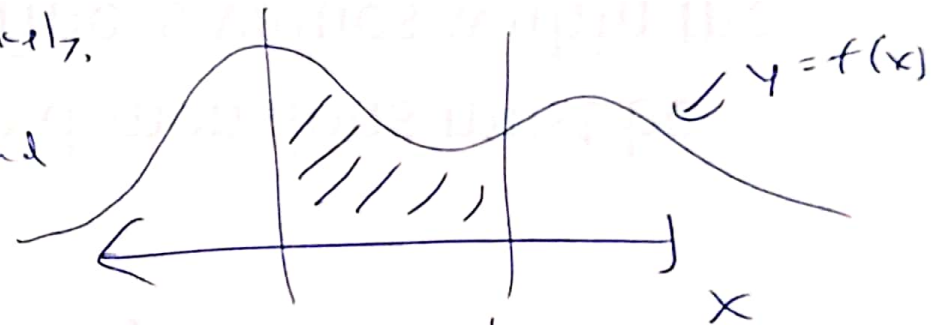② Kernel regression — Allows us to estimate

the [expected value] of y as a function

of x w/o assuming a particular form for   ) ex: linear velocity 引用

the regression curve. (It doesn't have to

be linear, quadratic, or any other

particular form.)

# The Probability Density Function pdf

Recall: The pdf $f(x)$ is the function that, for a continuous distribution, indicates which values are likely and which are unlikely. (It integrates to 1) and it completely determines the distribution, as does the distribution function.

$y = f(x)$

$P(a < X < b)$

$$= \int_a^b f(x)\, dx$$

One estimate: Create a histogram with equal-sized bins, and standardize it so that it integrates to 1. ✳

Some problems with this: Not smooth, sensitive to the choice of bins for the histogram.

# Parametric Density Estimation 有多法

⟹ **Q:** How would we estimate a (probability density) within a particular parametric family?

**A:** We would estimate the parameters & then use the corresponding density as our estimate.

For example, given $X_1, \ldots, X_n \overset{iid}{\sim} N(\mu, \sigma^2)$, we could compute $\bar{X}$ & $s$ (sufficient statistics) and then use the $N(\bar{x}, s^2)$ pdf as our estimate. (Such estimators would be consistent as long as the parameter estimates are consistent.) (for familiar families of distributions)
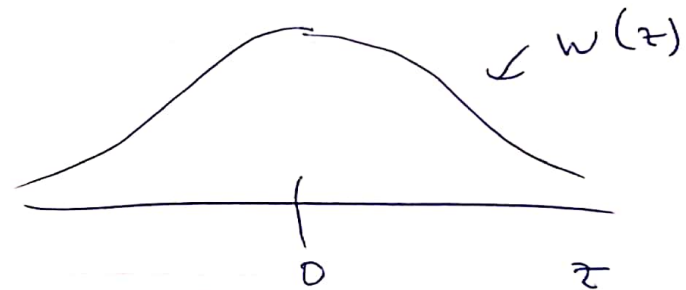
# Kernel Density Estimation ← nonparametric! 非参数

Let $X_1, X_2, \ldots, X_n$ be the data, a random sample from the ⁽ᶜᵀˢ⁾ distribution with pdf $f(x)$.

⟹ Let $w(z)$ be a __kernel__, a ① __symmetric__ density function ② with mean 0 ③ and standard deviation 1.

We then estimate $f(x)$ with

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{\Delta} w\left(\frac{x - X_i}{\Delta}\right),$$



↙ $w(z)$

One possible kernel.

note! where $\Delta$ is the bandwidth

Note: The book uses $h$ for the bandwidth, but $h$ looks too much like $n$ to work well in these notes.

∴ 省师 已望把 $\Delta$ 位为 bandwidths

# The Bandwidth △

If $w(z)$ has mean 0 and standard deviation 1, then $\frac{1}{\triangle} w\left(\frac{z}{\triangle}\right)$ has mean 0 and st. dev. △.
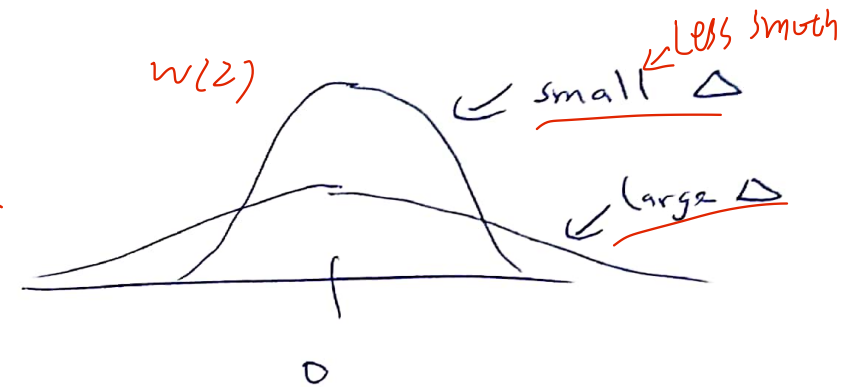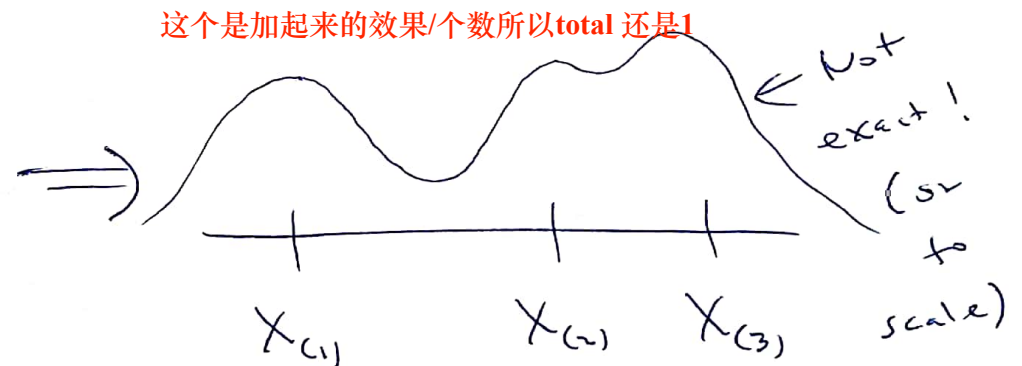
Big △ $\Rightarrow$ more spread out.

$w(z)$ ← small △ ← less smooth

← large △

0

→ Kernal density estimator

## Interpreting $\hat{f}(x)$:

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

We average together $n$ pdfs, one centered at each data value.

假设只有3个 data

这个是加起来的效果/个数所以total 还是1

← Not exact! (so to scale)

$X_{(1)}$ $X_{(2)}$ $X_{(3)}$ $\Rightarrow$ $X_{(1)}$ $X_{(2)}$ $X_{(3)}$

# Properties of $\hat{f}(x)$: <span style="color:red">kernal density</span>

① $\hat{f}(x)$ is a valid density. In particular, it is non-negative and integrates to 1.

Q: Why? A: Note that $\frac{1}{\Delta} w\left(\frac{x-x_i}{\Delta}\right)$ integrates to 1 for each $i$.

② $\hat{f}(x)$ is ⟨smooth⟩ if $w(\cdot)$ is smooth. ← $w$ need not be smooth, though!

③ If we let $\Delta \to 0$ sufficiently slowly as $n \to \infty$, then $\hat{f}(x)$ is a ⟨consistent⟩ estimator for $f(x)$ at a particular point.

# Choosing the Bandwidth

Bigger $\Delta$ leads to more averaging, smaller $\Delta$ to a more Jagged estimate.

Research indicates that $\Delta$ should go like $\frac{1}{n^{1/5}}$ to minimize the MSE for $\hat{E}(x)$.

One simple choice: (too simple?)

$$\Delta = \frac{1.06 \, s}{n^{1/5}} \quad \leftarrow \text{sample st. dev.}$$
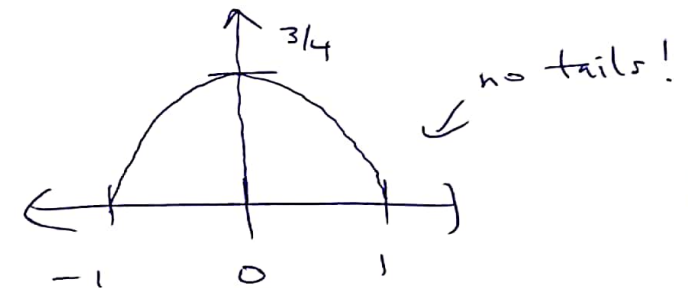
More sophisticated: We could choose $\Delta$ using leave-one-out cross-validation. We could even let $\Delta$ be different for different $x$ values.

# Some Kernel Choices

① **Gaussian:** Here the standard normal pdf is $w(z)$.

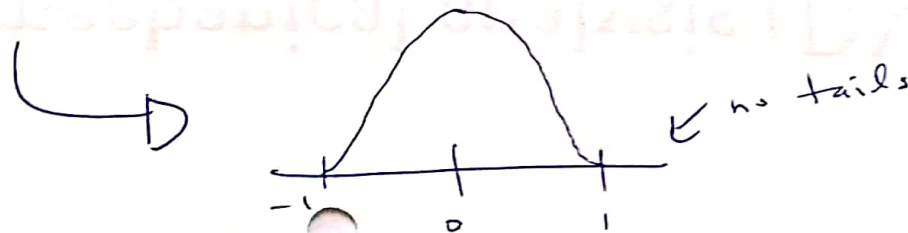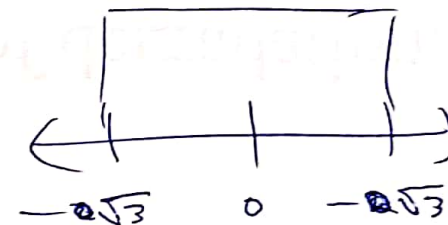② **Epanechnikov:** The kernel $w(z)$ is a standardized version

of the function $K(x) = \dfrac{3(1-x^2)}{4}$, $-1 \leq x \leq 1$.

<span style="color:red">绝对误差<br>相对最大<br>好很多</span>

asymptotically optimal in one sense

③ **Rectangular:** Here $w(z)$ is the

Uniform$(-\sqrt{3}, \sqrt{3})$ pdf.

④ **Biweight:** Here $w(z)$ is a

standardized version of

$K(x) = \dfrac{15}{16}(1-x^2)^2$, $-1 \leq x \leq 1$.

# Baseball Data Example

Try out the density function in R w/ different kernels and different bandwidths.

**Q:** What happens when the bandwidth becomes either very large or very small?

large：越smooth
small:越锯齿状

**Q:** Do you see any limitations for the kernel estimators that we are discussing?

In particular, do the estimators have any undesirable properties?

$\Rightarrow$ Kernel Regression ⓘ

Recall: The regression of $y$ on $x$ is the function $E[Y|X=x]$, the conditional expected value.

{ If the relationship is linear, then we get a regression line, but the relationship need not be linear.

?
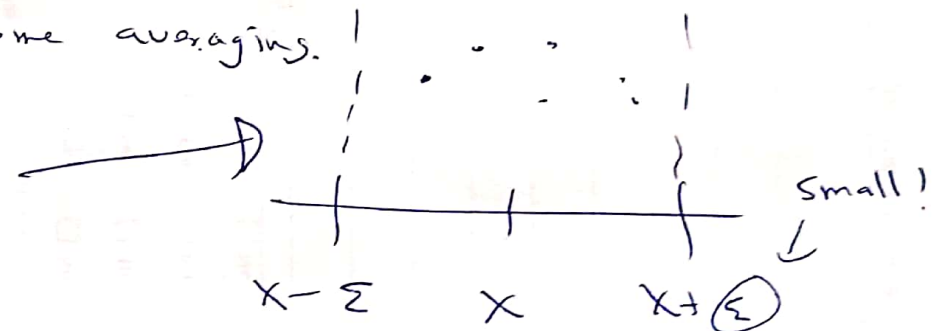
One can add higher-order terms like $x^2$, $x^3$, etc., but that also gives only a limited amount of flexibility.

One idea: If $E[Y|X=x]$ is continuous then it is roughly constant over a small interval of $x$ values.

This suggests doing some averaging.

Average these $y$ values to estimate $E[Y|X=x]$.

$\longrightarrow$

$x-\varepsilon$   $x$   $x+\varepsilon$

Small!

# Kernel Regression    (II)

One problem with this crude averaging is that $\widehat{E[Y|X=x]}$ won't be continuous in $x$ <u>at places where</u> one point <u>drops out of</u> the interval or joins the interval.

就是说中间断了一下

⟹ <u>Better Idea:</u> Use <u>smoothly - changing</u> <u>weights</u> based on a kernel, say $w(z)$.

Kernel regression estimates $\widehat{E[Y|X=x]} = \dfrac{\sum\limits_{i=1}^{n} Y_i \, w\!\left(\dfrac{x-x_i}{\Delta}\right)}{\sum\limits_{i=1}^{n} w\!\left(\dfrac{x-x_u}{\Delta}\right)}.$

Weight for $Y_i$ in the average

where $\Delta > 0$ is the bandwidth.

# Looking at the Weights

$$\downarrow^{x_1} \quad \downarrow^{x_2} \quad \downarrow^{x_3}$$

Suppose that we have X values $1, 3, 8$ and that we use a Gaussian kernel with $\Delta = 2$.

Using R, plot the weights for $y_1, y_2, + y_3$ as a function of x for x between $-5$ and $15$.

Q: What will happen to the estimates as x becomes very small or very large?

## Baseball Example:

Use kernel regression to estimate $E(Y|X=x)$ when $x =$ at bats and $y =$ batting average (success rate).

Try out different choices of the bandwidth.

Does it appear that the relationship is linear or not?

# Note on Bandwidth in R

In  density,  bw  is the  standard deviation of the kernel to be used.

In  ksmooth,  bw is set up  so that the quartiles of the kernel are $\pm 0.25$ bw.  This means that the standard deviation $\sigma$ satisfies,

$$bw \cong 2.7 \, \sigma.$$

Why?
$$0.25 \, bw \cong 0.674 \, \sigma$$
$$\implies bw \cong 2.7 \, \sigma.$$

# Tests for Contingency Tables

Recall: The chi-square test is used to test for association in a two-way table.

Test statistic: $X^2 = \mathcal{E} \dfrac{(Obs - Exp)^2}{Exp}$. If the

cells

expected counts are all large, then $X^2$ is

approximately distributed $X^2((r-1)\cdot(c-1))$, where

$r = \#$ of rows and $c = \#$ of columns.

Q: How large must the expected counts be?

A: Cochran suggests that all be at least 1.0, with no

more than 20% less than 5.0.

# Tests for Contingency Tables (II)

Q: What do we do if the expected counts are too small for the chi-square approximation to work well?

A: We can do an exact permutation-based version of the test.

Another possibility: Combining cells which is likely to lead to loss of power. It is also not always obvious which cells should be combined.

We condition on the row + column totals and find the permutation distribution of the test statistic.

The possibilities are not all equally likely!

# An Example

For the table $\begin{pmatrix} 4 & 0 \\ 1 & 4 \end{pmatrix}$ list out all the possible tables & find their probabilities under the null hypothesis of independence between the variables.

Solution:

$$\begin{pmatrix} 4 & 0 \\ 1 & 4 \end{pmatrix} \begin{matrix} 4 \\ 5 \end{matrix} \qquad \begin{pmatrix} 3 & 1 \\ 2 & 3 \end{pmatrix} \qquad \begin{pmatrix} 2 & 2 \\ 3 & 2 \end{pmatrix} \qquad \begin{pmatrix} 1 & 3 \\ 4 & 1 \end{pmatrix} \qquad \begin{pmatrix} 0 & 4 \\ 5 & 0 \end{pmatrix}$$

$$\frac{\binom{4}{4}\binom{5}{1}}{\binom{9}{5}} \qquad \frac{\binom{4}{3}\binom{5}{2}}{\binom{9}{5}} \qquad \frac{\binom{4}{2}\binom{5}{3}}{\binom{9}{5}} \qquad \frac{\binom{4}{1}\binom{5}{4}}{\binom{9}{5}} \qquad \frac{\binom{4}{0}\binom{5}{5}}{\binom{9}{5}}$$

$$= \frac{5}{126} \qquad = \frac{40}{126} \qquad = \frac{60}{126} \qquad = \frac{20}{126} \qquad = \frac{1}{126}$$

<span style="color:red">所有之和 =1</span>

Now find $X^2$ for each table, and do an upper-tailed test (like the chi-square test).

# Another Example

List out all the possibilities and probabilities for

the table $\begin{pmatrix} 2 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \begin{matrix} 3 \\ 2 \end{matrix}$

$$\begin{matrix} 3 & 1 & 1 \end{matrix}$$

Solution:

$$\begin{pmatrix} 3 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix} \qquad \begin{pmatrix} 2 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \qquad \begin{pmatrix} 2 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} \qquad \begin{pmatrix} 1 & 1 & 1 \\ 2 & 0 & 0 \end{pmatrix}$$

$$\frac{\binom{3}{3}\binom{2}{0}\binom{2}{1}}{\binom{5}{3}\binom{2}{1}} \qquad \frac{\binom{3}{2}\binom{2}{1}\binom{1}{1}\binom{1}{0}}{\binom{5}{3}\binom{2}{1}} \qquad \frac{\binom{3}{2}\binom{2}{1}\binom{1}{1}}{\binom{5}{3}\binom{2}{1}} \qquad \frac{\binom{3}{1}\binom{2}{2}\binom{2}{1}}{\binom{5}{3}\binom{2}{1}}$$

$$= \frac{2}{20} \qquad\qquad = \frac{6}{20} \qquad\qquad = \frac{6}{20} \qquad\qquad = \frac{6}{20}$$

These must sum to 1!

# Fisher's Exact Test

If the contingency table is 2 × 2, it is also possible to test independence (or equality of proportions) using Fisher's exact test.

The listing is done exactly as for the chi-square test, and (the test statistic is the value in the (1,1) cell or some equivalent test statistic.)

Here the test can be one-sided, unlike for the chi-square test.

Ex: In a game against South Florida in 2010, Corey Fisher made 1 three-point shot and missed 3, while Corey Stokes made 3 and missed 3. Find the p-values for ⓐ a two-sided test and ⓑ a one-sided test where the alternative is that Stokes succeeds more often.

# Fisher's Exact Test Example ①

Observed table:

$$\begin{array}{c} & \text{made} \quad \text{missed} \\ \text{Fisher} & \begin{pmatrix} 1 & 3 \\ 3 & 3 \end{pmatrix} \begin{matrix} 4 \\ 6 \end{matrix} \\ \text{Stokes} & \quad 4 \quad\quad 6 \end{array}$$

Possibilities:

$$\begin{pmatrix} 0 & 4 \\ 4 & 2 \end{pmatrix} \quad \begin{pmatrix} 1 & 3 \\ 3 & 3 \end{pmatrix} \quad \begin{pmatrix} 2 & 2 \\ 2 & 4 \end{pmatrix} \quad \begin{pmatrix} 3 & 1 \\ 1 & 5 \end{pmatrix} \quad \begin{pmatrix} 4 & 0 \\ 0 & 6 \end{pmatrix}$$

$$\frac{\binom{4}{0}\binom{6}{4}}{\binom{10}{4}} \quad \frac{\binom{4}{1}\binom{6}{3}}{\binom{10}{4}} \quad \frac{\binom{4}{2}\binom{6}{2}}{\binom{10}{4}} \quad \frac{\binom{4}{3}\binom{6}{1}}{\binom{10}{4}} \quad \frac{\binom{4}{4}\binom{6}{0}}{\binom{10}{4}}$$

$$210 \longrightarrow \binom{10}{4}$$

$$= \frac{15}{210} \qquad = \frac{80}{210} \qquad = \frac{90}{210} \qquad = \frac{24}{210} \qquad = \frac{1}{210}$$

Probability of a table as likely or ~~more~~ less likely than the observed table

ⓐ add all but $\frac{90}{210}$ to get $\frac{15+80+24+1}{210} \cong \boxed{0.57}$

ⓑ add the observed prob. and those to the left to get $\frac{80+15}{210} \cong \boxed{0.45}$

# Fisher's Exact Test Example (II)

Verify our results for this example in R.

Note: The way we computed the two-tailed p-value is valid in general, but in symmetric situations, it reduces to doubling the shorter of the two one-tailed p-values.