

$$\alpha_i = \mu_{i.} - \mu_{..}$$

$$\beta_{j(i)} = \mu_{ij} - \mu_{i.}$$

$$\sum_{i=1}^a \alpha_i = 0 \quad \sum_{j=1}^b \beta_{j(i)} = 0 \quad \text{for } i=1, \dots, a$$

ANOVA

partitioning the SS (balanced design)

$$y_{ijk} = \bar{y}_{...} + (\bar{y}_{i..} - \bar{y}_{...}) + (\bar{y}_{.j.} - \bar{y}_{i.}) + (y_{ijk} - \bar{y}_{ij.})$$

$$\begin{aligned} SSTO &= \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{...})^2 \\ &= \sum_i \sum_j \sum_k (\bar{y}_{i..} - \bar{y}_{...})^2 + \sum_i \sum_j \sum_k (\bar{y}_{.j.} - \bar{y}_{i.})^2 + \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{ij.})^2 \\ &= bn \sum_i (\bar{y}_{i..} - \bar{y}_{...})^2 + n \sum_i \sum_j (\bar{y}_{.j.} - \bar{y}_{i.})^2 + \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{ij.})^2 \\ &= SS_A + SS_B(A) + SSE \end{aligned}$$

Note in two-way ANOVA

$$\begin{aligned} SSTO &= bn \sum_i (\bar{y}_{i..} - \bar{y}_{...})^2 + an \sum_j (\bar{y}_{.j.} - \bar{y}_{...})^2 + n \sum_i \sum_j (\bar{y}_{ij.} - \bar{y}_{i.} - \bar{y}_{.j.} + \bar{y}_{...})^2 \\ &\quad + \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{ij.})^2 \end{aligned}$$

$$SS_{B(A)} = SS_B + SS_{AB}$$

Under Normality, all SS_s/σ^2 are indep.

Computationally.

$$SS_{TO} = \sum_i \sum_j \sum_k y_{ijk}^2 - \frac{y_{...}^2}{abn} \quad y_{...} = \sum_i \sum_j \sum_k y_{ijk}$$

$$SS_A = \frac{1}{bn} \sum_i y_{i..}^2 - \frac{y_{...}^2}{abn}$$

$$SS_{B(A)} = \frac{1}{n} \sum_j \sum_k y_{.jk}^2 - \frac{1}{bn} \sum_i y_{i..}^2$$

$$SSE = \sum_i \sum_j \sum_k y_{ijk}^2 - \frac{1}{n} \sum_i y_{i..}^2$$

Anova table.

Source	SS	df	MS
A	SS_A	$a-1$	
B(A)	$SS_{B(A)}$	$(b-1) + (a-1)(b-1)$ $= a(b-1)$	
Error	SSE	$ab(n-1)$	
total	SS_T		

• For nested design.

For any bracketed subscripts in the model, place a 1 under those subscripts that are inside the brackets.

• If A and B fixed.

	F		R	
	a	b	n	k
α_i	0	b	n	
$\beta_{j(i)}$	1	0	n	

$$E MS$$

$$E MSE = \sigma^2$$

$$E MSA = \sigma^2 + b n \sum_{i=1}^a \frac{\alpha_i^2}{a-1}$$

$$E MS_{B(A)} = \sigma^2 + n \frac{\sum_i \sum_j \beta_{j(i)}^2}{a(b-1)}$$

$$F\text{-test } A: \frac{MSA}{MSE}, \quad B(A): \frac{MS_{B(A)}}{MSE}$$

• If A Fixed, B random.

	R		R
α_i	a	b	Σ
	0	b	n
$\beta_{j(c)}$	1	1	n

F test

A:

$$\frac{MSA}{MSB(A)}$$

$$B(A): \frac{MSB(A)}{MSE}$$

$$\hat{\sigma}_\beta^2 = \frac{MSB(A) - MSE}{n - 1}$$

$$E(MSA) = \sigma^2 + b n \sigma_\alpha^2 + n \sigma_\beta^2$$

$$E(MSB(A)) = \sigma^2 + n \sigma_\beta^2$$

$$E(MSE) = \sigma^2$$

F-test A:

$$\frac{MSA}{MSB(A)}$$

B(A):

$$\frac{MSB(A)}{MSE}$$

$$\hat{\sigma}_\beta^2 = \hat{\sigma}_\alpha^2 =$$

• If A, B random

	R		R
α_i	1	b	n
$\beta_{j(c)}$	1	1	n

Traditional factorial design (Chap 29)

(136)

2^k factorial design

- k factors, each factor has two levels (often labeled + & -)
- very useful for preliminary analysis.
- can "remove" unimportant factors.
- all interactions.

$$k=2$$

general 2-factor factorial model

$$\mu_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$$
$$1 + (a-1) + (b-1) + (a-1)(b-1)$$

4 parameters.

$$\mu, \alpha_1 = -\alpha_2, \beta_1 = -\beta_2$$

$$(\alpha\beta)_{11} + (\alpha\beta)_{12} = 0$$

$$(\alpha\beta)_{21} + (\alpha\beta)_{22} = 0$$

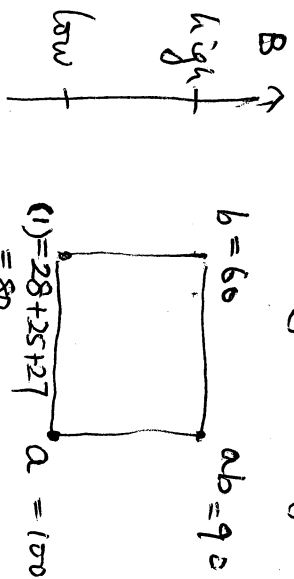
$$(\alpha\beta)_{11} + (\alpha\beta)_{21} = 0$$

$$(\alpha\beta)_{11} = (\alpha\beta)_{22}$$

Design for 2^2 factorial

- Label the levels of factors A and B using + and -
- There are 4 experimental combinations labeled.

Symbol	Factor			I	II	III
(1)	A	B				
a	+	-	A low B low	28	25	27
b	-	+	A high B low	36	32	32
ab	+	+	A low B high	18	19	23
			A high B high	31	30	29



- Can express combination in terms of model parameters.

$$ab = (+, +) = EY = \mu + \alpha_2 + \beta_2 + (\alpha\beta)_{22}$$

$$a = (+, -)$$

$$EY = \mu + \alpha_2 + \beta_1 + (\alpha\beta)_{21}$$

$$b = (-, +)$$

$$EY = \mu + \alpha_1 + \beta_2 + (\alpha\beta)_{12}$$

$$(1) = (-, -)$$

$$EY = \mu + \alpha_1 + \beta_1 + (\alpha\beta)_{11}$$

ab

$$E y = \mu + \alpha_2 + \beta_2 + (\alpha\beta)_{22}$$

(138)

a

$$E y = \mu + \alpha_2 - \beta_2 - (\alpha\beta)_{22}$$

b

$$E y = \mu + \alpha_2 + \beta_2 - (\alpha\beta)_{22}$$

(4)

$$E y = \mu - \alpha_2 - \beta_2 + (\alpha\beta)_{22}$$

Estimate parameters (4)

$$\hat{\mu} = \frac{ab + a + b + (1)}{4n}$$

n: replications
of each combination

$$\hat{\alpha}_2 = \frac{1}{2} \left(\frac{ab+a}{2n} - \frac{b+(1)}{2n} \right) = \frac{ab+a-b-(1)}{4n}$$

$$\hat{\beta}_2 = \frac{1}{2} \left(\frac{ab+b}{2n} - \frac{a+(1)}{2n} \right) = \frac{ab-a+b-(1)}{4n}$$

$$(\hat{\alpha}\hat{\beta})_{22} = \frac{1}{2} \left(\frac{ab+(1)}{2n} - \frac{a+b}{2n} \right) = \frac{ab-a-b+(1)}{4n}$$

Nonparametric Regression (kernel smoothing spline smoothing)

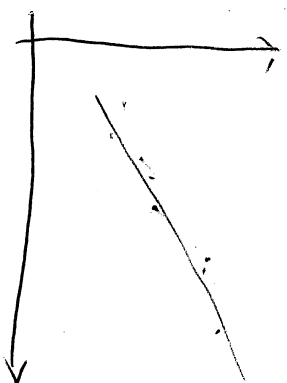
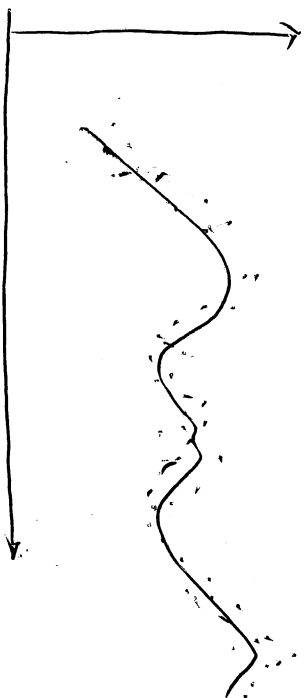
Regression problem: Given fixed (x_1, \dots, x_n) , we observe (y_1, \dots, y_n)

$$\text{where } y_i = m(x_i) + \varepsilon_i$$

\downarrow
iid

$$E(\varepsilon_i) = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2$$

m : unknown, the problem is to estimate m .



Assumption:

the nonparametric approach is to choose m from some smooth family of functions, "some degree of smoothness".

- advantage and shortcoming

pros: flexibility

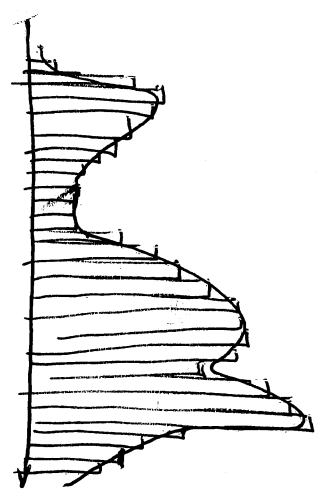
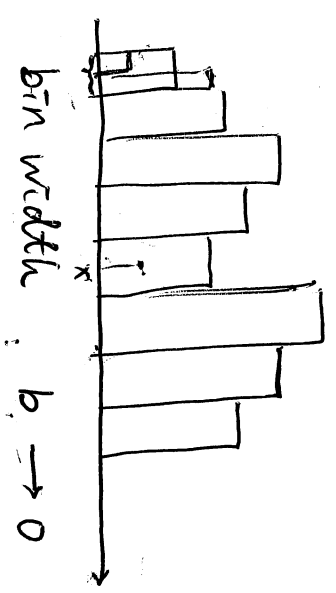
useful when little past experience is available

con: parametric approach is more effective if the model is correct.

- do not have a formulaic way to describe the relationship (graphically)

- kernel smoothing

histogram

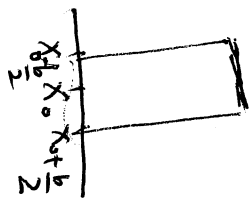


$f(x)$: prob. density func. at x .

(14)

$$\frac{\# \text{ of } x_i \text{ in the } B_n(x)}{n} \approx P(X \in B(x))$$

$$\int_a^b f(x) dx = P(a \leq X \leq b)$$



$$\int_{x_0 - \frac{b}{2}}^{x_0 + \frac{b}{2}} f(x) dx \approx f(x_0) \cdot b$$

$$\hat{f}(x_0, b) = \frac{P(X \in (x_0 - \frac{b}{2}, x_0 + \frac{b}{2}))}{b}$$

$$= \frac{\# \text{ of } x_i \text{ in } (x_0 - \frac{b}{2}, x_0 + \frac{b}{2})}{nb}$$

$$= \frac{\sum_{i=1}^n I[X_i \in (x_0 - \frac{b}{2}, x_0 + \frac{b}{2})]}{nb}$$

$$E \left[\sum_{i=1}^n I[X_i \in (x_0 - \frac{b}{2}, x_0 + \frac{b}{2})] \right] \approx nb f(x_0)$$

If $n \rightarrow \infty$, $nb f(x_0) \rightarrow \infty$

So $nb \rightarrow \infty$

③

$n \rightarrow \infty$
 $b \rightarrow 0$
①

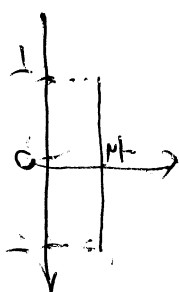
for example $b \sim \frac{1}{\sqrt{n}}$

(142)

h : bandwidth.

$$\hat{f}_h(x) = \frac{\sum_{i=1}^n I(X_i \in (x-h, x+h))}{n \cdot 2h}$$

$$= \frac{\sum_{i=1}^n \frac{1}{2h} I\left(\left|\frac{X_i - x}{h}\right| < 1\right)}{n}$$



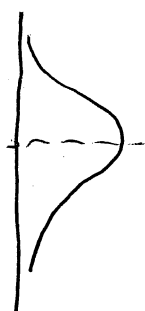
$$= \frac{\sum_{i=1}^n \frac{1}{n} K\left(\frac{X_i - x}{h}\right)}{n} \quad \leftarrow \text{here } K(u) = \frac{I(|u| < 1)}{2}$$

unif kernel

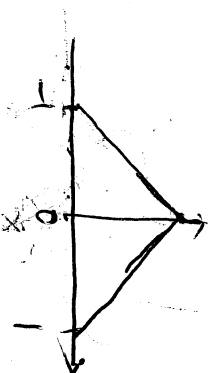
$$= \frac{\sum_{i=1}^n K_h(X_i - x)}{n} \quad \leftarrow \text{here } K_h(u) = \frac{1}{h} K\left(\frac{u}{h}\right)$$

Choice of kernel function

$K(u) = \varphi(u)$ standard normal



$K(u) = (1 - |u|) I(|u| < 1)$ triangular kernel



kernel function

i). $h \rightarrow 0$, $nh \rightarrow \infty$

ii) k is symmetric $\int_{-\infty}^{\infty} u k(u) du = 0$

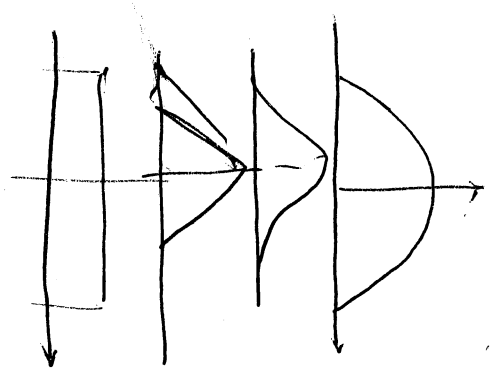
$$\int_{-\infty}^{\infty} k(u) du = 1$$

$$k(u) \geq 0 \quad , \quad \int u^2 k(u) du < \infty$$

kernel function selection.

Epanechnikov is the "best" kernel

Normal	.951
triangular	.986
uniform	.93



Density estimation

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

$$i) \hat{f}_h(x) \geq 0$$

$$ii) \int_R \hat{f}_h(x) dx = \frac{1}{n} \sum_{i=1}^n \int \frac{1}{h} K\left(\frac{x-X_i}{h}\right) dx = 1$$

$$MSE(\hat{f}_h(x)) = Bias^2 + Variance$$

$$\text{here } Bias = E[\hat{f}_h(x)] - f(x)$$

$$Variance = Var(\hat{f}_h(x))$$

$$Bias = E[\hat{f}_h(x)] - f(x) \rightarrow \{X_1, \dots, X_n\} \sim f(x)$$

$$= E\left[\frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) - f(x)\right]$$

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

$$\text{Let } v = \frac{u-x}{h} \rightarrow \int \frac{1}{h} K(-v) f(x+vh) h dv - f(x)$$

$$= \int K(v) f(x+vh) dv - f(x)$$

$$\stackrel{\text{Taylor}}{=} \int K(v) \left[f(x) + vh f'(x) + \frac{(vh)^2}{2} f''(x) + o(h^2) \right] dv - f(x)$$

$$= \cancel{f(x)} + f'(x) \underbrace{\int v K(v) dv}_0 + \frac{f''(x) h^2}{2} \int v^2 K(v) dv + o(h^2) - \cancel{f(x)}$$

$$= \frac{h^2 f''(x)}{2} \mu_2(K) + o(h^2)$$

$$\downarrow$$

$$\mu_m(K) = \int_{-\infty}^{\infty} K(v) v^m dv$$

