

Lecture 4

Multivariate visualization

- identify structure in data (e.g. existence of clusters of observations), outliers, etc.
- look at informative 2D scatterplots

Different approaches

- ① Find "interesting" 2D projections of the data.

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \quad A = (\underline{a}_1, \underline{a}_2) \quad \text{Typically } \underline{a}_1, \underline{a}_2 \text{ orthogonal: } \underline{a}_1^T \underline{a}_2 = 0$$

$$Y = XA = \begin{pmatrix} y_{11} & y_{12} \\ \vdots & \vdots \\ y_{n1} & y_{n2} \end{pmatrix}$$

Plot $\{y_{11}\}$ vs. $\{y_{12}\}$

- $\underline{a}_1, \underline{a}_2$ should be chosen so that $\{x_i^T \underline{a}_1\} + \{x_i^T \underline{a}_2\}$ are "interesting"
- ↖ identify groups, outliers etc.

- ② Grand Tour Method

Asimov (1980s)

- do an exhaustive tour of all 2D projections $t = \text{time parameters}$

$$A(t) = (\underline{a}_1(t), \underline{a}_2(t))$$

$$Y(t) = XA(t)$$



$n \times 2 \Rightarrow$ scatterplot movie

- how to define $\underline{a}_1(t)$ & $\underline{a}_2(t)$ as t varies

Simple visualizations

Scatterplot matrix: Look at all $\binom{P}{2} = \frac{P(P-1)}{2}$ pairwise scatterplots of variables.

In R: pairs-displays scatterplots in a matrix

- more informative version of a correlation matrix
- non-linear pairwise dependencies revealed

Downside: Miss out on higher dimensional structure.

Mathematically - we're looking at projections $\underline{a} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ b \\ \vdots \\ 0 \end{pmatrix} \rightarrow$ i^{th} element

How do we find "interesting" projections?

- how do we define "interesting"? Random vector $\underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \sim f(\underline{x}) \rightarrow$ joint dist'n

Sps that for any α the dist'n of $\alpha^T X$ is the same up to location & scale.
 i.e. for some $\mu(\alpha)$, $\sigma(\alpha)$, the dist'n of $\frac{\alpha^T X - \mu(\alpha)}{\sigma(\alpha)}$ is the same for all α

- suggests that all projections are uninteresting

Multivariate normal has the property: all one dimensional projections are normal!

"Def'n": Non-interesting projection = projections that are normal.

Look for non-normal projections

Projection pursuit: interesting

- try to find non-normal projections of the data.

How to do this? \rightarrow tractable

- need indices of non-normality

- $I(\alpha^T x_1, \dots, \alpha^T x_n)$ is large if $\{\alpha^T x_i\}$ is non-normal

\downarrow
Index

- try to find α to maximize I .

Example: Kurtosis

- for a random variable X with mean μ & variance σ^2 we define

$$\boxed{K(X) = \frac{E[(X-\mu)^4]}{\sigma^4}}$$

$$K(ax+b) = K(X)$$

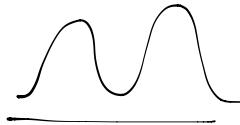
If $X \sim N(\mu, \sigma^2)$ then $K(X) = 3$

- short/light-tailed dist'n $K < 3$

- long/heavy-tailed dist'n $K > 3$

- multi-modal dist'n $K > 3$

modes? or model?



Estimating Kurtosis

data x_1, \dots, x_n with mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and covariance matrix $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$

$$\therefore \hat{K}(\alpha) = \frac{\frac{1}{n} \sum_{i=1}^n (\alpha^T x_i - \alpha^T \bar{x})^4}{(\alpha^T S \alpha)^2}$$

- Find α to maximize (or minimize) $\hat{K}(\alpha)$

- Alternative (and equivalently) we can maximize (or minimize) $\frac{1}{n} \sum_{i=1}^n (\alpha^T x_i - \alpha^T \bar{x})^4$ subject to $\alpha^T S \alpha = 1$

$$\frac{1}{n} \sum_{i=1}^n (\alpha^T x_i - \alpha^T \bar{x})^4 + \lambda [\alpha^T S \alpha - 1]$$

$$\frac{\partial}{\partial \alpha} : \frac{4}{n} \sum_{i=1}^n (\alpha^T x_i - \alpha^T \bar{x})^3 (x_i - \bar{x}) + 2\lambda S \alpha$$

$$\alpha = \frac{-4}{2\lambda n} \sum_{i=1}^n (\alpha^T x_i - \alpha^T \bar{x})^3 S^{-1} (x_i - \bar{x})$$

Example: Exam marks data 5 exam marks
 88 students

- maximized kurtosis = 5.26

$$\alpha = \begin{pmatrix} 0.18 \\ 0.76 \\ -0.61 \\ -0.09 \\ -0.02 \end{pmatrix}$$

Isolate one observation as an outlier

$$\tilde{x}_{81} = \begin{pmatrix} 3 \\ 9 \\ 51 \\ 47 \\ 40 \end{pmatrix} \quad \tilde{\bar{x}} = \begin{pmatrix} 39.0 \\ 30.6 \\ 50.6 \\ 46.7 \\ 42.3 \end{pmatrix}$$

- also look at $(\tilde{x}_i - \tilde{\bar{x}})^T S^{-1} (\tilde{x}_i - \tilde{\bar{x}})$ ← compare to $\chi^2(S)$ quantiles

- minimized Kurtosis = 2.06

$$\underline{\alpha} = \begin{pmatrix} -0.45 \\ 0.84 \\ 0.30 \\ 0.02 \\ -0.05 \end{pmatrix}$$

From 1 to 2 dimensions

- maximize sample kurtosis $\hat{K}(\underline{\alpha})$ at $\underline{\alpha} = \underline{\alpha}_1$

- now maximize $K(\underline{\alpha})$ over $\underline{\alpha}$ with $\underline{\alpha}^T \underline{\alpha}_1 = 0$

Maximize $\frac{1}{n} \sum_{i=1}^n (\underline{\alpha}^T \underline{x}_i - \underline{\alpha}^T \tilde{\bar{x}})^4$ subject to $\underline{\alpha}^T S \underline{\alpha} = 1$ and $\underline{\alpha}^T \underline{\alpha}_1 = 0$
 $\Rightarrow \underline{\alpha}_2$

Plot $\underline{\alpha}_1^T \underline{x}_i$ vs $\underline{\alpha}_2^T \underline{x}_i$ for $i=1, \dots, n$

Example: Exam mark data

- plots and R code on Blackboard

$$\underline{\alpha}_1 = \begin{pmatrix} -0.45 \\ -0.84 \\ 0.30 \\ 0.02 \\ -0.05 \end{pmatrix} \quad \underline{\alpha}_2 = \begin{pmatrix} 0.40 \\ -0.17 \\ 0.02 \\ -0.82 \\ 0.36 \end{pmatrix}$$

Plot of $\underline{\alpha}_1^T \underline{x}_i$ vs $\underline{\alpha}_2^T \underline{x}_i$ again reveals: \underline{x}_{81} as unusual

Density Estimation

Given univariate or bivariate data

$$x_1, \dots, x_n \quad \underline{x}_1, \dots, \underline{x}_n$$

- want to estimate underlying density f .

Kernel density estimation (1 dimension)

- symmetric kernel function $K(x)$ ← density

- $K(-x) = K(x)$
- $K(x) \geq 0$
- $\int_{-\infty}^{\infty} K(x) dx = 1$

- for some bandwidth parameter h , define $\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$

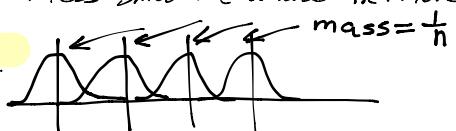
- R function: density

- Bandwidth h controls smoothness of $\hat{f}_h(x)$

- increase $h \Rightarrow$ smoother density estimate

- decrease $h \Rightarrow$ less smooth estimate i.e. more modes

Motivation



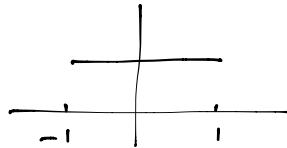
- spread out the probability mass of $\frac{1}{n}$ at each point
- at each x_i , add up contributions from each x_i
- around x_i , prob max($=\frac{1}{n}$) is redist'd according to the density function

$$\begin{aligned} & \frac{1}{nh} K\left(\frac{x-x_i}{h}\right) \\ \text{- note that } & \int_{-\infty}^{\infty} \frac{1}{nh} K\left(\frac{x-x_i}{h}\right) dx, \quad y = \frac{x-x_i}{h} \quad dx = h dy \\ & = \frac{1}{n} \int_{-\infty}^{\infty} K(y) dy = \frac{1}{n} \end{aligned}$$

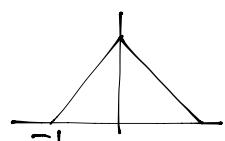
- add contributions from each x_i : $\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$

Types of kernels

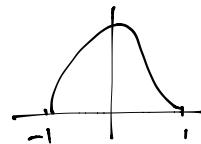
① Rectangular kernel $K(x) = \frac{1}{2}$ for $|x| \leq 1$



② Triangular kernel $K(x) = 1 - |x|$ for $|x| \leq 1$



③ Gaussian Kernel $K(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$



④ Epanechnikov kernel $K(x) = \frac{3}{4}(1-x^2)$ for $|x| \leq 1$

Notes

① Choice of kernel is much less important than bandwidth choice

- rectangular kernel should be avoided

- bias/variance tradeoff in choice of h

- as $h \uparrow$, bias \uparrow , variance \downarrow

$h \downarrow$, bias \downarrow , variance \uparrow

- given the true f , can derive optimal values of h

- automatic bandwidth selection in R is good.

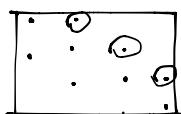
② R function density, scales the kernels to have variance 1

$$\text{i.e. } \int_{-\infty}^{\infty} x^2 K(x) dx = 1$$

- makes choice of bandwidth for different kernels

easy \rightarrow same h gives very similar estimates for different kernels.

From 1 to 2 dimensions (briefly)



- kernel density estimation more complicated in 2 dimensions