

WHAT STATISTICIANS DO

THIS IS AN ATTEMPT TO EXPLAIN WHAT STATISTICS IS BY LOOKING AT WHAT STATISTICIANS DO, AS OPPOSED TO WHAT YOU OR THEY SHOULD OR SHOULDN'T DO.

1. IDEALLY, STATISTICIANS SHOULD START WITH A DESIGN PHASE: HOW LARGE A SAMPLE IS NEEDED, WHERE SHOULD OBSERVATIONS BE TAKEN, WILL WE HAVE THE RIGHT DATA FOR THE ANSWERS WE NEED. THIS IS A PRACTICAL RANITY, THOUGH ^{EXPERIMENTAL} DESIGN IS A BOOMING BUSINESS, OFTEN MORE IN THEORY THAN PRACTICE. PARAPHRASINGLY, MUCH PROGRESS IN DESIGN CENTERED AROUND THE FOLLOWING PROBLEM: WE HAVE A SPACE WITH A DISTANCE, WE WANT TO PUT n POINTS DOWN SO THAT THE MAXIMUM OF THE DISTANCE TO ANY OTHER POINT IS MINIMIZED.
2. THE BULK OF STATISTICAL WORK IS ANALYSIS, AND THIS IS CRUCIAL TO ANSWERING THE DESIGN QUESTION. WHEN STATISTICIANS COMMUNICATE TO EACH OTHER, THE FOLLOWING TREE-LIKE DESCRIPTION IS OFTEN IMPLICIT.
3. DATA TYPE 1 IS THE DATA A SAMPLE FROM A WELL DEFINED LARGER POPULATION, IF SO WHAT KIND OF SAMPLE (SIMPLE, PROPORTIONAL TO A ZONE^{ZONE}...)... IF NOT, WE CALL THE DATA SET A SAMPLE OF CONVIVENCE OR A BATCH AND WORRY ABOUT THE VALUE OF OUR ANALYSIS WHEN APPLIED TO ANOTHER POPULATION. TODAY, THE DATA MAY OFTEN ARISES AS THE OUTPUT OF A COMPLEX COMPUTER SIMULATION; AT ANOTHER EXTREME, THE DATA MAY BE A COMPLETE ENUMERATION OF A FINITE ^{POPULATION}.
4. DATA TYPE 2 THE DATA MAY BE UNIVARIATE LENGTHS OF A SAMPLE OF FISH, BIVARIATE (x_i, y_i) HEIGHT AND WEIGHT FOR EXAMPLE, OR MULTIVARIATE. IT MAY BE A TIME SERIES (STOCK PRICES) OR DATA ALONG THE

2

GENOME, IT MAY BE FROM A SPATIAL PROCESS. (LOCATIONS OF INCIDENCE OF CANCER ON A MAP).

4. DATA TYPE 3 WE OFTEN HAVE DIFFERENT ANALYTICAL TOOLS FOR CONTINUOUS DATA VERSUS DISCRETE DATA (CONTINGENCY TABLES). IF EACH IS PRESENT WE TALK ABOUT MIXED DATA. NOW ADAYS WE ALSO SEE DATA IN THE FORM OF CURVES (e.g. GAIT ANALYSIS, GROWTH CURVES OF CHILDREN) AND DATA IN THE FORM OF RANKINGS OR PHYLOGENETIC TREES!

CLASSICAL.

5. STATISTICAL TASKS THERE ARE THREE CLASSICAL STATISTICAL TASKS: ESTIMATION (WHERE IS THIS BATCH OF DATA CENTERED), TESTING (DOES THE SUGGESTED MODEL FIT THE DATA) AND REGRESSION GIVEN SOME TRAINING DATA $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ FIND A FUNCTION $f(x)$ WHICH ~~DOES~~ MAKES $f(x_i)$ CLOSE TO y_i . THIS FUNCTION MAY BE USED FOR PREDICTION OR FOR UNDERSTANDING (DOES INCREASING MONEY TO SCHOOLS AFFECT STUDENT PERFORMANCE).

6. MODERN STATISTICAL TASKS THE PURPOSE OF INVESTIGATION MAY BE JUST TO GET A HUMINELY INTERPRETABLE DESCRIPTION OF THE DATA AND ITS RELATIONS. THIS IS OFTEN CALLED EXPLORATORY DATA ANALYSIS (EDA), AND IS OFTEN CENTERED AROUND GRAPHICAL DISPLAYS, WAYS OF VISUALIZING IN HIGH DIMENSIONS. A RECENT HYBRID WHICH IS WIDELY PRACTICED IS MODEL BUILDING; THIS IS THE ITERATIVE CYCLE OF FITTING A PRELIMINARY MODEL (BY ESTIMATING A CURVE ON SOME PARAMETER) CHECKING ITS GOODNESS OF FIT (BY TESTING OR GRAPHICAL ANALYSIS OF RESIDUALS) CHANGING THE MODEL TO INCORPORATE HIGHER ORDER APPROXIMATIONS, AND SO ON.

7. TYPES OF MODELS THOSE MAY BE PARAMETRIC (e.g. FIT A GAMMA DISTRIBUTION $f(x|\alpha, \sigma^2) = \frac{1}{\sigma} \alpha^{-1} x^{\alpha-1} e^{-\frac{x}{\sigma}}$ ON $\{\alpha, \sigma\}$) NOW PARAMETRIC

(FIND AN ESTIMATE FOR THE BEST FITTING CURVE WITH $\int f''^2 \leq 10$) OR SEMI-PARAMETRIC (ESTIMATE THE BEST MONOTONE FUNCTION OR THE BEST SYMMETRIC DENSITY (THE CENTER OF SYMMETRY BEING THE PARAMETER)).

8. INFERENTIAL MODE. FREQUENTIST ANALYSIS Tries TO CHOOSE PROCEDURES WHICH WILL DO WELL IF MANY INSTANCES OF THE PROBLEM UNDER INDEPENDENT REALIZATIONS. FOR EXAMPLE, THE AVERAGE VALUE MINIMIZES THE AVERAGE SQUARED ERROR ABOUT THE TRUE MEAN (FOR UNIVARIATE DATA). Bayesian Analysis INVOLVES PUTTING A PRIOR DISTRIBUTION OVER ANY UNKNOWN PARAMETERS AND CHOOSES PARAMETER ESTIMATES TO MINIMIZE THE BAYES RISK. THUS, IN A BAYESIAN SETUP, ONE HAS A FAMILY OF PROBABILITIES p_θ ON AND A PRIOR DISTRIBUTION $\pi(\theta)$; TO ESTIMATE θ_n , IT IS USUAL TO SPECIFY A LOSS FUNCTION $L(\theta, \hat{\theta}(x))$. WE CHOOSE AN ESTIMATE TO MINIMIZE $\int L(\theta, \hat{\theta}(x)) p_\theta d\theta$. THIS IS OFTEN THE MEAN OR MODE OF THE POSTERIOR DISTRIBUTION $\pi(\theta|x)$ GIVEN BY BAYES THEOREM.

I WANT TO MENTION TWO OTHER INFERENTIAL MODES: DETERMINISTIC -- WILL A SMALL CHANGE IN THE DATA EFFECT MY CONCLUSION? SUSAN HOLMES THINKS THE BOOTSTRAP ANALYSIS IS MOSTLY USED FROM THIS POINT OF VIEW. FINALLY, THERE IS THE BAYESIAN VIEW: THEOREM TELLS US THAT A MEDIUM TO LARGE AMOUNT OF DATA, THE CONCLUSIONS DRAWN FROM BAYESIAN OR FREQUENTIST ANALYSIS WILL AGREE TO GOOD APPROXIMATION. ONE MAY USE A BAYES PROCEDURE TO GET SOMETHING GOING WITHOUT TAKING THE PRIOR AS A SERIOUS QUANTIFICATION OF OPINION KNOWLEDGE. ALTERNATIVELY, MANY BAYESIANS CHOOSE THE PARAMETERS IN THEIR PRIORS FROM THE AVAILABLE DATA (EMPIRICAL BAYES).

9. THE ABOVE CATEGORIES ABOVE GIVE A TREE LIKE DESCRIPTION TO STATISTICS; FOLLOWING A PATH GIVES VARIOUS AREAS WITHIN THE FIELD: BAYESIAN NON-PARAMETRIC ESTIMATION OF UNIVARIATE TIME SERIES DATA,

4

10. THERE ARE ~~many~~^{SOME} considerations LEFT OUT. ONE OF THE BIGGEST AREAS OF STATISTICAL DEVELOPMENT TWENTY YEARS AGO WAS ROBUST ANALYSIS DEVELOPING AUTOMATED WAYS TO DEAL WITH OUTLIERS. AN IMPORTANT SUB-AREA IS MISSING DATA AND SURVIVAL ANALYSIS; WHAT TO DO IF SOME DATA ARE MISSING. MATHEMATICAL STATISTICS TRIES TO MAKE MATHEMATICALLY PLEASING CHOICES OF GOOD PROCEDURES ON A GIVEN BRANCH OF THE TREE. IT SUPPOSES GIVEN A FAMILY $p_\theta(x)$ AND A LOSS FUNCTION L_C