# Lecture Stat 302
## Introduction to Probability - Slides 7

AD

Feb. 2010

# Simpson's Paradox: Sex Bias in Graduate Admissions?

- The University of California at Berkeley was sued for bias against women who had applied for admission to graduate schools there. The admission figures for the fall of 1973 showed that men applying were more likely than women to be admitted, and the difference was so large that it was unlikely to be due to chance.

- The evidence looks really compelling.

|  | Applicants | % admitted |
|---|---|---|
| Men | 8442 | 44% |
| Women | 4321 | 35% |

- If we consider the event $A =$ "Admitted", $B =$ "Be a woman", $B^c =$ "Be a man", we have

$$P\left(A|\,B\right) \leq P\left(A|\,B^c\right).$$

- Was Berkeley really biased?

# Simpson's Paradox: Sex Bias in Graduate Admissions?

- Let us examine more carefully the data by examining admissions in some representative departments

| Depart. | Men % in | Women % in |
|---------|----------|------------|
| 1 | 63 | 68 |
| 2 | 33 | 35 |
| 3 | 6 | 7 |

- Whatever the department to which they apply, women have an higher probability of getting in than men; i.e. we have for any event $C_i = \{\text{apply to department } i\}$

$$P\left(A|\, B \cap C_i\right) \geq P\left(A|\, B^c \cap C_i\right)$$

and still overall

$$P\left(A|\, B\right) \leq P\left(A|\, B^c\right)$$

# Simpson's Paradox: Sex Bias in Graduate Admissions?

- Let us examine even more carefully the data

| Depart. | Men Appli. | % in | Women Appli. | % in |
|---------|------------|------|--------------|------|
| 1 | 560 | 63 | 25 | 68 |
| 2 | 417 | 33 | 375 | 35 |
| 3 | 272 | 6 | 341 | 7 |

# Simpson's Paradox: Sex Bias in Graduate Admissions?

- The mathematical reason is that $P(A|B \cap C_i) \geq P(A|B^c \cap C_i)$ does NOT indeed imply $P(A|B) \geq P(A|B^c)$ as

$$P(A|B) = \sum_{i=1}^{n} P(A|B \cap C_i) P(C_i|B),$$

$$P(A|B^c) = \sum_{i=1}^{n} P(A|B^c \cap C_i) P(C_i|B^c),$$

It highly depends on the proba of applying to department $i$ given you are a woman/man .

- The reason is that women tended to apply to competitive departments with low rates of admission, whereas men tended to apply to less-competitive departments with high rates of admission.

# Proof of expression for conditional proba

- This is very similar to proof of $P(A) = \sum_{i=1}^{n} P(A| C_i) P(C_i)$.
- We have
$$P(A| B) = P(A \cap (\cup_{i=1}^{n} C_i)| B)$$
and $A \cap (\cup_{i=1}^{n} C_i) = \cup_{i=1}^{n} (A \cap C_i)$ where $(A \cap C_i) \cap (A \cap C_j) = \varnothing$ for $i \neq j$ so
$$P(A| B) = \sum_{i=1}^{n} P(A \cap C_i| B)$$
where by the definition of conditional proba
$$\begin{aligned} P(A \cap C_i| B) &= \frac{P(A \cap C_i \cap B)}{P(B)} \\ &= \frac{P(A \cap C_i \cap B)}{P(C_i \cap B)} \frac{P(C_i \cap B)}{P(B)} \\ &= P(A| B \cap C_i) P(C_i| B). \end{aligned}$$

# Simpson's Paradox: A "Smoking" Example

- In 1972-1994 a survey of the electoral roll, largely concerned with smoking habits and survival rates was carried out in Wichkham, a mixed urban and rural district near Newcastle upon Tyne, in the UK. Twenty years later, a follow-up study was conducted.
- Relationship between smoking habits and 20-year survival in 401 women aged 55-74.

|       | Smoker | Non-Smoker |
|-------|--------|------------|
| Dead  | 80     | 141        |
| Alive | 71     | 109        |

- Introducing $A=$"die", $B=$"smoking", we have

$$P(A|B) = \frac{80}{80+71} = 0.53,$$
$$P(A|B^c) = \frac{141}{141+109} = 0.56.$$

- Does smoking really help living longer?

# Simpson's Paradox: A "Smoking" Example

- Now lets look more precisely at the data. We have an additional variable related to the age.
- We have

| Age 55-64 | Smoker | Non-Smoker | | Age 65-74 | Smoker | Non-Smoker |
|-----------|--------|------------|---|-----------|--------|------------|
| Dead | 51=44% | 40=33% | \| | Dead | 29=80% | 101=78% |
| Alive | 64=56% | 81=67% | | Alive | 7=20% | 28=22% |

- For both age class, the survival rate is smaller for the non-smokers. Mathematically, this means that if we introduce $C = \{$age between 55-64$\}$ and then $C^c = \{$age between 65-74$\}$ then we have as expected

$$P\left(A|\, B \cap C\right) \;\geq\; P\left(A|\, B^c \cap C\right),$$
$$P\left(A|\, B \cap C^c\right) \;\geq\; P\left(A|\, B^c \cap C^c\right)$$

but still we have

$$P\left(A|\, B\right) \leq P\left(A|\, B^c\right).$$

# Simpson's Paradox: A "Smoking" Example

- This should not be a surprise as

$$P(A|B) = P(A|B \cap C)P(C|B) + P(A|B \cap C^c)P(C^c|B),$$
$$P(A|B^c) = P(A|B^c \cap C)P(C|B^c) + P(A|C^c \cap B^c)P(C^c|B^c)$$

so the results depend on $P(C|B), P(C^c|B), P(C|B^c), P(C^c|B^c)$ i.e. that is the proba of being in a given age given smoking or not!

- The "paradox" occurs as a higher proportion of non-smokers studied belong to older age groups, in which survival rates for both smokers and non-smokers are significantly lower compared to younger age groups.

- In other words, most of the smokers have died off before reaching the older age classes and so the higher number of deaths (in absolute numbers) for the non-smokers in the older age classes has obscured the result.

# Random Variables

- In many scenarios, we are interested in a function of the outcome as opposed to the actual outcome; e.g. we are interested in the sum of two dice and not in the separate values of each die or, when we flip a coin, we want to know the number of tails.

- Real-valued functions defined on the sample space are *random variables*; e.g. number of gold medals by Canadian athletes at Olympics, your score at the SAT test etc.

- **Example**: We are tossing 3 times a fair coin. If we call $X$ the number of heads obtained, then it is a random variables such that

$$P\{X=0\} = P(\{T,T,T\}) = \left(\frac{1}{2}\right)^3 = \frac{1}{8},$$
$$P\{X=1\} = P(\{T,T,H\},\{T,H,T\},\{H,T,T\}) = \frac{3}{8},$$
$$P\{X=2\} = P(\{T,H,H\},\{H,T,H\},\{H,H,T\}) = \frac{3}{8},$$
$$P\{X=3\} = P(\{H,H,H\}) = \frac{1}{8}.$$

## Examples

- **Coin Toss**: For a coin toss, the possible events are heads or tails. The number of heads appearing in one fair coin toss can be described using the following random variable:

$$X = \left\{ \begin{array}{ll} 0 & \text{if head} \\ 1 & \text{if tail} \end{array} \right. \quad \text{and } P\left(X=0\right) = P\left(X=1\right) = \frac{1}{2}$$

Remark: Note that we could define completely arbitrary random variable such as $X = 10$ if heads and $X = \pi$ is tail. This is of no "practical" interest but possible conceptually.

- **Rolling a fair die**: A random variable can also be used to describe the process of rolling a fair dice and the possible outcomes. The most obvious representation is to take the set $\left\{1, 2, 3, 4, 5, 6\right\}$ as the sample space, defining the random variable $X$ as the number rolled. In this case,

$$X = i \text{ if a } i \text{ is rolled and } P\left(X=i\right) = \frac{1}{6} \text{ for } i = 1, ..., 6.$$

# Example

- Assume you have trials consisting of sitting an exam. You have proba $p$ of passing the exam. If you fail, you sit the exam again. These trials are considered independent as you never bother revising the material. We denote $X$ the number of exams you need to sit before passing the subject.

- We have

$$
\begin{aligned}
P\left(X=1\right) &= P\left(Pass\right) = p, \\
P\left(X=2\right) &= P\left(Fail, Pass\}\right) = (1-p)\,p, \\
&\quad .... \\
P\left(X=n\right) &= P(\underbrace{Fail, ..., Fail}_{n-1 \text{ times}}, Pass) = (1-p)^{n-1}\,p.
\end{aligned}
$$

- You can check that

$$
\sum_{k=1}^{\infty} P\left(X=k\right) = \sum_{k=1}^{\infty} (1-p)^{k-1}\,p = 1.
$$

# Example: Roulette Wheel

- The pockets of the roulette wheel are numbered from 1 to 36. In number ranges from 1 to 10 and 19 to 28, odd numbers are red and even are black. In ranges from 11 to 18 and 29 to 36, odd numbers are black and even are red. There is two green pockets numbered 0 and 00.
- *Red/Black*: If you bet 1\$ on red, the payout is 1\$ if it is red. Let $X_1$ denotes the payout

$$P\left(X_1 = 1\right) = P\left(\text{red}\right) = \frac{18}{18+18+2} = 0.4737$$
$$P\left(X_1 = -1\right) = 1 - P\left(X = 1\right) = 0.5263$$

as red$=\{1, 5, 7, 9, 12, 14, 16, 18, 19, 21, 23, 25, 27, 30, 32, 34, 36\}$.

- *Straight up*: you bet 1\$ on any pocket 0,00,1,2,...,36. The payout is is 35\$ if you are right. Let $X_2$ denotes the payout

$$P\left(X_2 = 35\right) = P\left(\text{outcome is the pocket bet on}\right) = \frac{1}{38},$$

$$P\left(X_2 = -1\right) = P\left(\text{outcome is not on the pocket bet on}\right) = \frac{37}{38}.$$

# Example: Roulette Wheel

- *Split bet*: you bet 1\$ on 11 or 14. The payout is 17\$ if you are right. Let $X_3$ denotes the payout

$$P\left(X_3 = 17\right) = P\left(\text{outcome is 11 or 14}\right) = \frac{2}{38},$$
$$P\left(X_3 = -1\right) = P\left(\text{outcome is neither 11 nor 14}\right) = \frac{36}{38}.$$

- *Street bet*: you 1\$ on 1,2,3 or 4,5,6 or... or 19,20,21. Say you bet on 19,20,21. The payout is 11\$ if you are right. Let $X_4$ denotes the payout

$$P\left(X_4 = 11\right) = P\left(\text{outcome is 19, 20 or 21}\right) = \frac{3}{38},$$
$$P\left(X_4 = -1\right) = P\left(\text{outcome is not 19, 20 or 21}\right) = \frac{35}{38}.$$

# Discrete Random Variables

- A random variable (r.v.) $X$ which can take at most a countable number of possible values is said to be *discrete*.

- The *probability mass function* of $X$ is denoted by

$$p(a) = P(X = a).$$

- If $X$ can take the values $\{x_1, x_2, ...\}$ then

$$p(x_i) \geq 0 \text{ and } \sum_{i=1}^{\infty} p(x_i) = 1$$

- *Cumulative distribution function* is

$$F(a) = P(X \leq a) = \sum_{x:x \leq a} p(x).$$

# Example: Poisson distribution

- Assume $X$ is a discrete r.v. taking values $\{0, 1, 2, ...\}$ where $p(i) = C\lambda^i/i!$ with $\lambda > 0$. 1) Find the expression of $C$, 2) Compute $P(X = 0)$ and $P(X \geq 2)$

- We have

$$\sum_{i=0}^{\infty} p(i) = C \sum_{i=1}^{\infty} \lambda^i/i! = 1$$

but $e^\lambda = \sum_{i=1}^{\infty} \lambda^i/i!$ so

$$C = e^{-\lambda}.$$

- It follows that $P(X = 0) = e^{-\lambda}\frac{\lambda^0}{0!} = e^{-\lambda}$ and

$$
\begin{aligned}
P(X \geq 2) &= 1 - P(X < 2) = 1 - P(X = 0) - P(X = 1) \\
&= 1 - e^{-\lambda} - \lambda e^{-\lambda}.
\end{aligned}
$$