

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Binomial Distribution & Sampling

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## What is a Random Variable?

Formally,

$$R : \mathcal{S} \rightarrow \mathbb{R}$$

Sample space

(usually)

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Independent Variables

Random variables  $R, S$

are **independent** iff

$$[R = a], [S = b]$$

are independent *events*

for all numbers  $a, b$ .

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Independent Variables

Alternative version:

$$\Pr\{R = a \text{ and } S = b\} = \Pr\{R = a\} \cdot \Pr\{S = b\}.$$

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Mutually Independent RV's

**Mutual** Independence of  
random vars  $A_1, A_2, \dots, A_n$ :

$$\Pr\{A_1=a_1 \text{ and } A_2=a_2 \text{ and } \dots A_n=a_n\} = \Pr\{A_1=a_1\} \cdot \Pr\{A_2=a_2\} \cdots \Pr\{A_n=a_n\}.$$

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Independent Variables

**k-wise** Independence:  
any  $k$  of the variables are  
mutually independent  
(2-wise = **pairwise**)

6	9	13	7
12	10	5	
3	7	4	14
15	8	11	2

## Independent Variables

Pairwise Independence  
sufficient for major  
applications (in later lecture)  
which is useful since pairwise  
holds in important cases where  
mutual does not.

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.7

6	9	13	7
12	10	5	
3	7	4	14
15	8	11	2

## Density & Distribution

The **Probability Density Function**  
of random variable  $R$ ,

$$\text{PDF}_R(a) ::= \Pr\{R=a\}$$

**Cumulative Distribution Function** of  $R$ ,

$$\text{CDF}_R(a) ::= \Pr\{R \leq a\}$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.8

6	9	13	7
12	10	5	
3	7	4	14
15	8	11	2

## Indicator Variables

**Indicator variable** for event  $A$ :

$$I_A ::= \begin{cases} 1 & \text{if } A \text{ occurs,} \\ 0 & \text{if } \overline{A} \text{ occurs.} \end{cases}$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.11

6	9	13	7
12	10	5	
3	7	4	14
15	8	11	2

## Distributions

*Example:*

$H_i ::=$  indicator for a head on  
the  $i$ th coin flip.

Coin may be *biased*:

$$\Pr\{H_i = 1\} = p \neq 1/2$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.12

6	9	13	7
12	10	5	
3	7	4	14
15	8	11	2

## Binomial Distribution

$H_{n,p} ::=$  # heads in  $n$  mutually  
independent flips of a  $p$ -biased  
coin.

$$H_{n,p} = H_1 + H_2 + \dots + H_n$$

Probability space: the  $2^n$   
sequences of  $n$  H's and T's.

$$\Pr\{Q\} ::= p^{\#H's \text{ in } Q} \cdot (1-p)^{\#T's}$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.14

6	9	13	7
12	10	5	
3	7	4	14
15	8	11	2

## Binomial Distribution

$\Pr\{k \text{ Heads}\} =$   
(# $k$  head seqs)  
 $\cdot \Pr\{\text{seq with } k \text{ H's}\}$

$$\text{PDF}_{H_{n,p}}(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.16

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Polling & Sampling

Estimate % contaminated fish in Charles River?



*Procedure:* catch  $n$  fish, test each, use %contaminated in catch as estimate of %contaminated in whole river

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.18

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Sampling Questions



Catch 100 fish; what is probability that estimate is within 10% of actual%?

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.19

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Model as Coin Tosses



$p$  ::= fraction contaminated in river  
Fish tested: coin toss with bias  $p$ .

Catching  $n$  fish: tossing  $n$  coins

$A$  ::= fraction contaminated in the sample of 100

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.20

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Polling using Binomial PDF

$A = \# \text{"heads"}/100$   
within 10% of  $p$ ?

$$\Pr\{|A - p| \leq 0.1\} = \Pr\{|H_{100,p} - 100p| \leq 10\}$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.21

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Polling using Binomial PDF

How do we bound this probability when we don't know  $p$ ?

**Lemma:**  $\Pr\{|H_{n,p} - 100p| \leq 10\}$  is min for  $p = 1/2$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.22

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Compute the exact probability

$$\Pr\{|A - p| \leq 0.1\} \geq$$

$$\Pr\{|H_{100,1/2} - 50| \leq 10\}$$

$$= \sum_{h=40}^{60} \binom{100}{h} 2^{-100} \geq 0.96$$

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.23

6	9	10	7
12	16	5	
3	5	4	14
15	8	11	2

## Confidence

We can be **96% confident** that our estimated fraction is within **0.1** of the actual fraction of contaminated fish in the whole river.

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.24

6	9	10	7
12	16	5	
3	5	4	14
15	8	11	2

## Sample size for better estimate

Suppose we want an estimate of the fraction that will be **4% ( $\pm 0.04$ )** accurate for **95%** of the time? Similar calculation implies need to sample **589** fish.

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.25

6	9	10	7
12	16	5	
3	5	4	14
15	8	11	2

## Confidence – **not** Probable Reality

Now suppose we sample **589** fish and discover **47** are contaminated. So we estimate  $p$  is **47/589**.

It's tempting to say

~~"the probability that~~

~~$p = 47/589 \pm 0.04$~~

~~is at least 95%"~~

**--Technically not correct!**

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.26

6	9	10	7
12	16	5	
3	5	4	14
15	8	11	2

## Confidence

$p$  is the actual fraction of bad fish in the river.

$p$  is **unknown**, but not a random variable!

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.27

6	9	10	7
12	16	5	
3	5	4	14
15	8	11	2

## Confidence

The possible outcomes of our *sampling procedure* is a random variable. We can say that "the **probability** that **our sample fraction** will be within  $\pm 0.04$  of the true fraction is at least **95%**"

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.28

6	9	10	7
12	16	5	
3	5	4	14
15	8	11	2

## Confidence

For simplicity we say that

$p = 47/589 \pm 0.04$  at the **95% confidence level**

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M.29

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Binomial Approximation

Numerical approximations  
for  $\text{PDF}_{H_{n,p}}(\alpha n)$ ,  
 $\text{CDF}_{H_{n,p}}(\alpha n)$ ,  
in Notes 13.

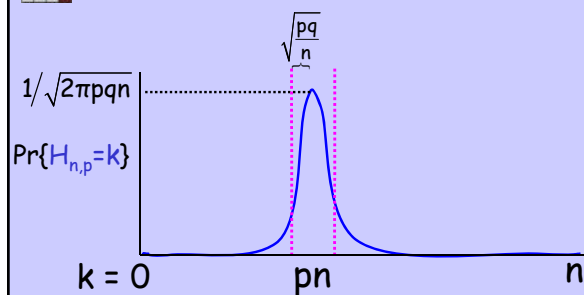
Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M-20

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Distribution of Heads



Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M-21

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Binomial Approximation

Messy formulas, but **easy to compute**.

*Exact* answers for  $n$  more than a few 1000 are impossible to compute  
(requires arithmetic on million-digit numbers)

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M-24

6	9	13	7
12	10	5	
3	4	8	14
15	2	11	1

## Team Problems

# Problems 1&2

Copyright © 2007, Albert R. Meyer. All rights reserved.

May 7, 2007

lec-13M-35