# DATA SIENCE

Lecture 4

# Lecture 4: Stochastic Thinking

# Newtonian Mechanics

- Every effect has a cause

- The world can be understood causally



1643 - 1727

# Copenhagen Doctrine

- Copenhagen Doctrine (Bohr and Heisenberg) of causal nondeterminism

  - At its most fundamental level, the behavior of the physical world cannot be predicted.

  - Fine to make statements of the form "x is highly likely to

  - occur," but not of the form "x is certain to occur."

# Stochastic Processes

- An ongoing process where the next state might depend on both the previous states and some random element

```python
def rollDie():
    """returns an int between 1 and 6"""


def rollDie():
    """returns a randomly chosen int
        between 1 and 6"""
```

# Implementing a Random Process

```python
import random

def rollDie():
    """returns a random int between 1 and 6"""
    return random.choice([1,2,3,4,5,6])

def testRoll(n = 10):
    result = ''
    for i in range(n):
        result = result + str(rollDie())
    print(result)
```

# Implementing a Random Process

```python
import random

def rollDie():
    """returns a random int between 1 and 6"""
    return random.choice([1,2,3,4,5,6])

def testRoll(n = 10):
    result = ''
    for i in range(n):
        result = result + str(rollDie())
    print(result)
```

# Probability Is About Counting

- Count the number of possible events

- Count the number of events that have the property of interest

- Divide one by the other

- Probability of 11111?
  - 11111, 11112, 11113, ..., 11121, 11122, ..., 66666
  - 1/(6**5)
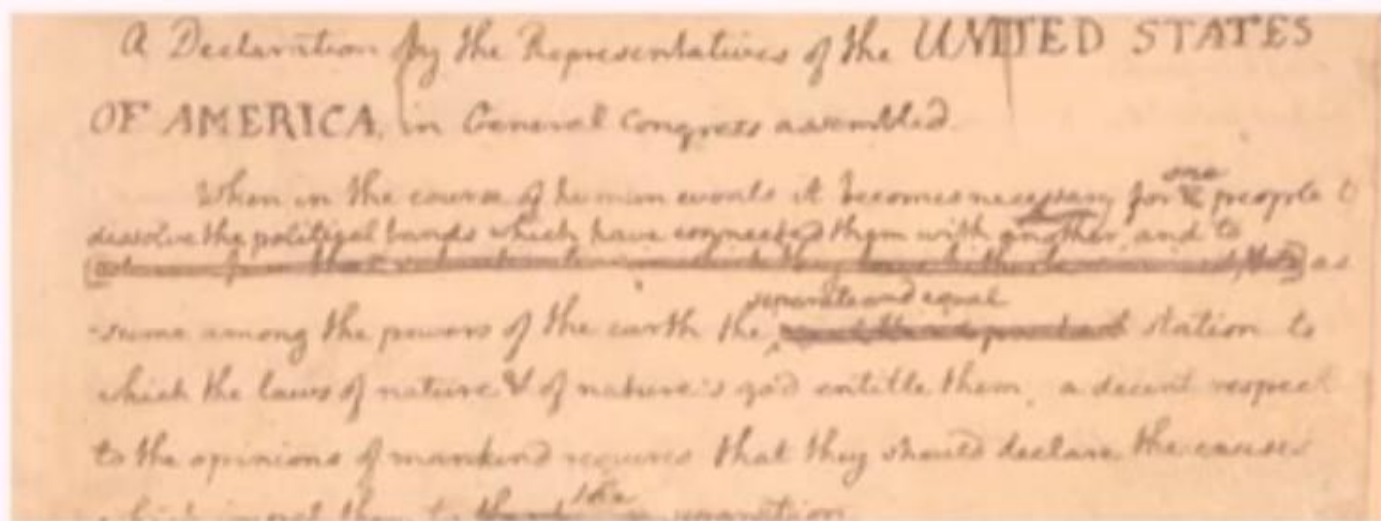
# Three Basic Facts About Probability

- Probabilities are always in the range 0 to 1. 0 if impossible, and 1 if guaranteed.

- If the probability of an event occurring is p, the probability of it not occurring must be $1 - p$

- When events are independent of each other, the probability of all of the events occurring is equal to a product of the probabilities of each of the events occurring.

$$A = .5 \qquad B \quad .4$$

$$A \& B = .2$$

# Independence

- Two events are **independent** if the outcome of one event has no influence on the outcome of the other
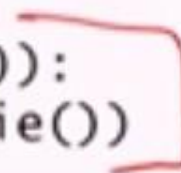
- Independence should not be taken for granted

# Will One of the Patriots and Broncos Lose?

- Patriots have winning percentage of 7/8, Broncos of 6/8

- Probability of both winning next Sunday is 7/8 * 6/8 = 42/64

- Probability of at least one losing is 1 – 42/64 = 22/64

# A Simulation of Die Rolling

```python
def runSim(goal, numTrials, txt):
    total = 0
    for i in range(numTrials):
        result = ''
        for j in range(len(goal)):
            result += str(rollDie())
        if result == goal:
            total += 1
    print('Actual probability of', txt, '=',
          round(1/(6**len(goal)), 8))
    estProbability = round(total/numTrials, 8)
    print('Estimated Probability of', txt, '=',
          round(estProbability, 8))

runSim('11111', 1000, '11111')
```

```
20
21 random.seed(0)
22
23 def runSim(goal, numTrials, txt):
24     total = 0
25     for i in range(numTrials):
26         result = ''
27         for j in range(len(goal)):
28             result += str(rollDie())
29         if result == goal:
30             total += 1
31     print('Actual probability of', txt, '=',
32         round(1/(6**len(goal)), 8))
33     estProbability = round(total/numTrials, 8)
34     print('Estimated Probability of', txt, '=',
35         round(estProbability, 8))
36
37 runSim('11111', 1000, '11111')
38
39 def sameDate(numPeople, numSame):
40     possibleDates = range(366)
41 #     possibleDates = 4*list(range(0, 57)) + [58]\
```

TypeError: cannot
instantiate ctype
'EVP_MD_CTX' of unknown
size


In [5]: runfile('C:/Users/
John/Dropbox (MIT)/current/
mit/Teaching/600/
Fall16/6.0002/lecture4/
lecture 4/lecture4.py',
wdir='C:/Users/John/Dropbox
(MIT)/current/mit/Teaching/
600/Fall16/6.0002/lecture4/
lecture 4')
Actual probability of 11111
= 0.0001286
Estimated Probability of
11111 = 0.0

In [6]:

# Output of Simulation

▪Actual probability = 0.0001286

▪Estimated Probability = 0.0

▪Actual probability = 0.0001286

▪Estimated Probability = 0.0

▪How did I know that this is what would get printed?

pseudo
random

seed

random.seed(0)

```python
20
21 random.seed(0)
22
23 def runSim(goal, numTrials, txt):
24     total = 0
25     for i in range(numTrials):
26         result = ''
27         for j in range(len(goal)):
28             result += str(rollDie())
29         if result == goal:
30             total += 1
31     print('Actual probability of', txt, '=',
32           round(1/(6**len(goal)), 8))
33     estProbability = round(total/numTrials, 8)
34     print('Estimated Probability of', txt, '=',
35           round(estProbability, 8))
36
37 runSim('11111', 1000000, '11111')
38
39 def sameDate(numPeople, numSame):
40     possibleDates = range(366)
41 #    possibleDates = 4*list(range(0, 57)) + [58]\
```

Console output (right panel):

```
600/Fall16/6.0002/lecture4/
lecture 4')
Actual probability of 11111
= 0.0001286
Estimated Probability of
11111 = 0.0

In [6]: runfile('C:/Users/
John/Dropbox (MIT)/current/
mit/Teaching/600/
Fall16/6.0002/lecture4/
lecture 4/lecture4.py',
wdir='C:/Users/John/Dropbox
(MIT)/current/mit/Teaching/
600/Fall16/6.0002/lecture4/
lecture 4')
Actual probability of 11111
= 0.0001286
Estimated Probability of
11111 = 0.000128

In [7]:
```

# The Birthday Problem

- What's the probability of at least two people in a group having the same birthday

- If there are 367 people in the group?

- What about smaller numbers?

- If we assume that each birthdate is equally likely
  - $1 - \dfrac{366!}{366^N * (366-N)!}$

- Without this assumption, VERY complicated

shoutkey.com/niece

# Approximating Using a Simulation

```python
def sameDate(numPeople, numSame):
    possibleDates = range(366)
    birthdays = [0]*366
    for p in range(numPeople):
        birthDate = random.choice(possibleDates)
        birthdays[birthDate] += 1
    return max(birthdays) >= numSame
```

# Approximating Using a Simulation

```python
def birthdayProb(numPeople, numSame, numTrials):
    numHits = 0
    for t in range(numTrials):
        if sameDate(numPeople, numSame):
            numHits += 1
    return numHits/numTrials

for numPeople in [10, 20, 40, 100]:
    print('For', numPeople,
          'est. prob. of a shared birthday is',
          birthdayProb(numPeople, 2, 10000))
numerator = math.factorial(366)
denom = (366**numPeople)*math.factorial(366-numPeople)
print('Actual prob. for N = 100 =',
      1 - numerator/denom)
```

Spyder (Python 3.5)

File  Edit  Search  Source  Run  Debug  Consoles  Projects  Tools  View  Help

Editor – C:\Users\John\Dropbox (MIT)\current\mit\Teaching\600\Fall16\6.0002\lecture4\lecture4.py

```
47          birthdays[birthDate] += 1
48      return max(birthdays) >= numSame
49
50 def birthdayProb(numPeople, numSame,
51     numHits = 0
52     for t in range(numTrials):
53         if sameDate(numPeople, numSa
54             numHits += 1
55     return numHits/numTrials
56
57 import math
58
59 for numPeople in [10, 20, 40, 100]:
60     print('For', numPeople,
61         'est. prob. of a shared bi
62         birthdayProb(numPeople, 2,
63     numerator = math.factorial(366)
64     denom = (366**numPeople)*math.fa
65     print('Actual prob. for N = 100
66         1 - numerator/denom)
67
```

IPython console

Console 1/A

```
current/mit/Teaching/600/Fall16/6.0002/lecture4/
lecture 4/lecture4.py', wdir='C:/Users/John/
Dropbox (MIT)/current/mit/Teaching/600/
Fall16/6.0002/lecture4/lecture 4')
Actual probability of 11111 = 0.0001286
Estimated Probability of 11111 = 0.000128

In [7]: runfile('C:/Users/John/Dropbox (MIT)/
current/mit/Teaching/600/Fall16/6.0002/lecture4/
lecture 4/lecture4.py', wdir='C:/Users/John/
Dropbox (MIT)/current/mit/Teaching/600/
Fall16/6.0002/lecture4/lecture 4')
For 10 est. prob. of a shared birthday is 0.1183
Actual prob. for N = 100 = 0.1166454118039999
For 20 est. prob. of a shared birthday is 0.4116
Actual prob. for N = 100 = 0.4105696370550831
For 40 est. prob. of a shared birthday is 0.8941
Actual prob. for N = 100 = 0.89054476188945
For 100 est. prob. of a shared birthday is 1.0
Actual prob. for N = 100 = 0.9999996784357714

In [8]:
```

Python console    History log    IPython console

Permissions: RW    End of lines: LF    Encoding: UTF-8    Line 59    Column 36    Memory 48 %

# Why 3 Is Much Harder Mathematically

- For 2 the complementary problem is "all birthdays distinct"

- For 3 people, the complementary problem is a complicated disjunct
  - All birthdays distinct or
  - One pair and rest distinct or
  - Two pairs and rest distinct or
  - ...

```
46          birthDate = random.choice(possibleDate
47          birthdays[birthDate] += 1
48      return max(birthdays) >= numSame
49
50  def birthdayProb(numPeople, numSame, numTrials
51      numHits = 0
52      for t in range(numTrials):
53          if sameDate(numPeople, numSame):
54              numHits += 1
55      return numHits/numTrials
56
57  import math
58
59  for numPeople in [10, 20, 40, 100]:
60      print('For', numPeople,
61          'est. prob. of a shared birthday is'
62          birthdayProb(numPeople, 3, 10000))
63      numerator = math.factorial(366)
64      denom = (366**numPeople)*math.factorial(36
65      print('Actual prob. for N = 100 =',
66          1 - numerator/denom)
67
```

```
Fall16/6.0002/lecture4/lecture 4/
lecture4.py', wdir='C:/Users/John/
Dropbox (MIT)/current/mit/Teaching/
600/Fall16/6.0002/lecture4/lecture 4')
For 10 est. prob. of a shared birthday
is 0.0011
Actual prob. for N = 100 =
0.1166454118039999
For 20 est. prob. of a shared birthday
is 0.0066
Actual prob. for N = 100 =
0.4105696370550831
For 40 est. prob. of a shared birthday
is 0.0651
Actual prob. for N = 100 =
0.89054476188945
For 100 est. prob. of a shared
birthday is 0.6359
Actual prob. for N = 100 =
0.9999996784357714

In [10]:
```

# Another Win for Simulation

- Adjusting analytic model a pain
- Adjusting simulation model easy

```
def sameDate(numPeople, numSame):
    possibleDates = 4*list(range(0, 57)) + [58]\
                          + 4*list(range(59, 366))\
                          + 4*list(range(180, 270))
    birthdays = [0]*366
    for p in range(numPeople):
        birthDate = random.choice(possibleDates)
        birthdays[birthDate] += 1
    return max(birthdays) >= numSame
```

```python
35              round(estProbability, 8))
36
37 #runSim('11111', 1000000, '11111')
38
39 def sameDate(numPeople, numSame):
40 #    possibleDates = range(366)
41     possibleDates = 4*list(range(0, 57)) + [58
42                     + 4*list(range(59, 366))\
43                     + 4*list(range(180, 270))
44     birthdays = [0]*366
45     for p in range(numPeople):
46         birthDate = random.choice(possibleDate
47         birthdays[birthDate] += 1
48     return max(birthdays) >= numSame
49
50 def birthdayProb(numPeople, numSame, numTrials
51     numHits = 0
52     for t in range(numTrials):
53         if sameDate(numPeople, numSame):
54             numHits += 1
55     return numHits/numTrials
56
```

Console output:

```
600/Fall16/6.0002/lecture4/lecture 4/
lecture4.py', wdir='C:/Users/John/
Dropbox (MIT)/current/mit/Teaching/
600/Fall16/6.0002/lecture4/lecture 4')
For 10 est. prob. of a shared birthday
is 0.0019
Actual prob. for N = 100 =
0.1166454118039999
For 20 est. prob. of a shared birthday
is 0.0097
Actual prob. for N = 100 =
0.4105696370550831
For 40 est. prob. of a shared birthday
is 0.0871
Actual prob. for N = 100 =
0.89054476188945
For 100 est. prob. of a shared
birthday is 0.7501
Actual prob. for N = 100 =
0.9999996784357714

In [11]:
```

# Simulation Models

- A description of computations that provide useful information about the possible behaviors of the system being modeled