

## Viva La Correlación!

- Say  $X$  and  $Y$  are arbitrary random variables
  - Correlation of  $X$  and  $Y$ , denoted  $\rho(X, Y)$ :
 
$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$
  - Note:  $-1 \leq \rho(X, Y) \leq 1$
  - Correlation measures linearity between  $X$  and  $Y$
  - $\rho(X, Y) = 1 \Rightarrow Y = aX + b$  where  $a = \sigma_Y/\sigma_X$
  - $\rho(X, Y) = -1 \Rightarrow Y = aX + b$  where  $a = -\sigma_Y/\sigma_X$
  - $\rho(X, Y) = 0 \Rightarrow$  absence of linear relationship
    - But,  $X$  and  $Y$  can still be related in some other way!
  - If  $\rho(X, Y) = 0$ , we say  $X$  and  $Y$  are “uncorrelated”
    - Note: Independence implies uncorrelated, but not vice versa!

## Fun with Indicator Variables

- Let  $I_A$  and  $I_B$  be indicators for events  $A$  and  $B$ 

$$I_A = \begin{cases} 1 & \text{if } A \text{ occurs} \\ 0 & \text{otherwise} \end{cases} \quad I_B = \begin{cases} 1 & \text{if } B \text{ occurs} \\ 0 & \text{otherwise} \end{cases}$$
- $E[I_A] = P(A)$ ,  $E[I_B] = P(B)$ ,  $E[I_A I_B] = P(AB)$
- $\text{Cov}(I_A, I_B) = E[I_A I_B] - E[I_A] E[I_B]$ 

$$= P(AB) - P(A)P(B)$$

$$= P(A | B)P(B) - P(A)P(B)$$

$$= P(B)[P(A | B) - P(A)]$$
- $\text{Cov}(I_A, I_B)$  determined by  $P(A | B) - P(A)$
- $P(A | B) > P(A) \Rightarrow \rho(I_A, I_B) > 0$
- $P(A | B) = P(A) \Rightarrow \rho(I_A, I_B) = 0$  (and  $\text{Cov}(I_A, I_B) = 0$ )
- $P(A | B) < P(A) \Rightarrow \rho(I_A, I_B) < 0$

## Can't Get Enough of that Multinomial

- Multinomial distribution
  - $n$  independent trials of experiment performed
  - Each trials results in one of  $m$  outcomes, with respective probabilities:  $p_1, p_2, \dots, p_m$  where  $\sum_{i=1}^m p_i = 1$
  - $X_i =$  number of trials with outcome  $i$
$$P(X_1 = c_1, X_2 = c_2, \dots, X_m = c_m) = \binom{n}{c_1, c_2, \dots, c_m} p_1^{c_1} p_2^{c_2} \dots p_m^{c_m}$$
  - E.g., Rolling 6-sided die multiple times and counting how many of each value  $\{1, 2, 3, 4, 5, 6\}$  we get
  - Would expect that  $X_i$  are negatively correlated
  - Let's see... when  $i \neq j$ , what is  $\text{Cov}(X_i, X_j)$ ?

## Covariance and the Multinomial

- Computing  $\text{Cov}(X_i, X_j)$ 
    - Indicator  $I_i(k) = 1$  if trial  $k$  has outcome  $i$ , 0 otherwise
 
$$E[I_i(k)] = p_i \quad X_i = \sum_{k=1}^n I_i(k) \quad X_j = \sum_{k=1}^n I_j(k)$$
    - $\text{Cov}(X_i, X_j) = \sum_{a=1}^n \sum_{b=1}^n \text{Cov}(I_i(b), I_j(a))$
    - When  $a \neq b$ , trial  $a$  and  $b$  independent:  $\text{Cov}(I_i(b), I_j(a)) = 0$
    - When  $a = b$ :  $\text{Cov}(I_i(b), I_j(a)) = E[I_i(a)I_j(a)] - E[I_i(a)]E[I_j(a)]$
    - Since trial  $a$  cannot have outcome  $i$  and  $j$ :  $E[I_i(a)I_j(a)] = 0$
- $$\text{Cov}(X_i, X_j) = \sum_{a=1}^n \sum_{b=1}^n \text{Cov}(I_i(b), I_j(a)) = \sum_{a=1}^n (-E[I_i(a)]E[I_j(a)])$$
- $$= \sum_{a=1}^n (-p_i p_j) = -n p_i p_j \Rightarrow X_i \text{ and } X_j \text{ negatively correlated}$$

## Multinomials All Around

- Multinomial distributions:
  - Count of strings hashed into buckets in hash table
  - Number of server requests across machines in cluster
  - Distribution of words/tokens in an email
  - Etc.
- When  $m$  (# outcomes) is large,  $p_i$  is small
  - For equally likely outcomes:  $p_i = 1/m$ 

$$\text{Cov}(X_i, X_j) = -n p_i p_j = -\frac{n}{m^2}$$
  - Large  $m \Rightarrow X_i$  and  $X_j$  very mildly negatively correlated
  - Poisson paradigm applicable

## Conditional Expectation

- $X$  and  $Y$  are jointly discrete random variables
  - Recall conditional PMF of  $X$  given  $Y = y$ :
 
$$p_{X|Y}(x|y) = P(X = x | Y = y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$
  - Define conditional expectation of  $X$  given  $Y = y$ :
 
$$E[X | Y = y] = \sum_x x P(X = x | Y = y) = \sum_x x p_{X|Y}(x|y)$$
  - Analogously, jointly continuous random variables:
 
$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} \quad E[X | Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx$$

## Rolling Dice

- Roll two 6-sided dice  $D_1$  and  $D_2$ 
  - $X = \text{value of } D_1 + D_2$      $Y = \text{value of } D_2$
- What is  $E[X | Y = 6]$ ?

$$E[X | Y = 6] = \sum_x xP(X = x | Y = 6)$$

$$= \left(\frac{1}{6}\right)(7 + 8 + 9 + 10 + 11 + 12) = \frac{57}{6} = 9.5$$

- Intuitively makes sense:  $6 + E[\text{value of } D_1] = 6 + 3.5$

## Hyper for the Hypergeometric

- $X$  and  $Y$  are independent random variables
  - $X \sim \text{Bin}(n, p)$      $Y \sim \text{Bin}(n, p)$
- What is  $E[X | X + Y = m]$ , where  $m \leq n$ ?
- Start by computing  $P(X = k | X + Y = m)$ :

$$P(X = k | X + Y = m) = \frac{P(X = k, X + Y = m)}{P(X + Y = m)} = \frac{P(X = k, Y = m - k)}{P(X + Y = m)} = \frac{P(X = k)P(Y = m - k)}{P(X + Y = m)}$$

$$= \frac{\binom{n}{k} p^k (1-p)^{n-k} \cdot \binom{n}{m-k} p^{m-k} (1-p)^{n-(m-k)}}{\binom{2n}{m} p^m (1-p)^{2n-m}} = \frac{\binom{n}{k} \binom{n}{m-k}}{\binom{2n}{m}}$$

- Hypergeometric:  $(X | X + Y = m) \sim \text{HypG}(m, 2n, n)$
- $E[X | X + Y = m] = nm/2n = m/2$     # total draws    total balls    white balls

## Properties of Conditional Expectation

- $X$  and  $Y$  are jointly distributed random variables

$$E[g(X) | Y = y] = \sum_x g(x) p_{X|Y}(x | y) \quad \text{or} \quad \int_{-\infty}^{\infty} g(x) f_{X|Y}(x | y) dx$$

- Expectation of conditional sum:

$$E\left[\sum_{i=1}^n X_i | Y = y\right] = \sum_{i=1}^n E[X_i | Y = y]$$

## Expectations of Conditional Expectations

- Define  $g(Y) = E[X | Y]$ 
  - $g(Y)$  is a random variable
  - For any  $Y = y$ ,  $g(Y) = E[X | Y = y]$ 
    - This is just function of  $Y$ , since we sum over all values of  $X$
  - What is  $E[E[X | Y]] = E[g(Y)]$ ? (Consider discrete case)

$$E[E[X | Y]] = \sum_y E[X | Y = y] P(Y = y)$$

$$= \sum_y \left[ \sum_x x P(X = x | Y = y) \right] P(Y = y)$$

$$= \sum_y \sum_x x P(X = x, Y = y) = \sum_x \sum_y x P(X = x, Y = y)$$

$$= \sum_x x P(X = x) = E[X] \quad (\text{Same for continuous})$$

## Analyzing Recursive Code

```
int Recurse() {
    int x = randomInt(1, 3); // Equally likely values
    if (x == 1) return 3;
    else if (x == 2) return (5 + Recurse());
    else return (7 + Recurse());
}
```

- Let  $Y = \text{value returned by Recurse}()$ . What is  $E[Y]$ ?

$$E[Y] = E[Y | X = 1]P(X = 1) + E[Y | X = 2]P(X = 2) + E[Y | X = 3]P(X = 3)$$

$$E[Y | X = 1] = 3$$

$$E[Y | X = 2] = E[5 + Y] = 5 + E[Y]$$

$$E[Y | X = 3] = E[7 + Y] = 7 + E[Y]$$

$$E[Y] = 3(1/3) + (5 + E[Y])(1/3) + (7 + E[Y])(1/3) = (1/3)(15 + 2E[Y])$$

$$E[Y] = 15$$

## Random Number of Random Variables

- Say you have a web site: [PimentoLoaf.com](http://PimentoLoaf.com)
  - $X = \text{Number of people/day visit your site. } X \sim N(50, 25)$
  - $Y_i = \text{Number of minutes spent by visitor } i. Y_i \sim \text{Poi}(8)$
  - $X$  and all  $Y_i$  are independent
  - Time spent by all visitors/day:  $W = \sum_{i=1}^X Y_i$ . What is  $E[W]$ ?

$$E[W] = E\left[\sum_{i=1}^X Y_i\right] = E\left[E\left[\sum_{i=1}^X Y_i | X\right]\right] = E[X \cdot E[Y_i]] = E[X]E[Y_i] = 50 \cdot 8$$

$$E\left[\sum_{i=1}^X Y_i | X = n\right] = \sum_{i=1}^n E[Y_i | X = n] = \sum_{i=1}^n E[Y_i] = nE[Y_i]$$

$$E\left[\sum_{i=1}^X Y_i | X\right] = X \cdot E[Y_i]$$



## Conditional Variance

- Recall definition:  $\text{Var}(X) = E[(X - E[X])^2]$ 
  - Define:  $\text{Var}(X | Y) = E[(X - E[X | Y])^2 | Y]$
- Derived:  $\text{Var}(X) = E[X^2] - (E[X])^2$ 
  - Can derive:  $\text{Var}(X | Y) = E[X^2 | Y] - (E[X | Y])^2$
- After a bit more math (in the book):
  - $\text{Var}(X) = E[\text{Var}(X | Y)] + \text{Var}(E[X | Y])$
  - Intuitively, let  $Y$  = true temperature,  $X$  = thermostat value
  - Variance in thermostat readings depends on:
    - Average variance in thermostat at different temperatures +
    - Variance in average thermostat value at different temperatures

## Making Predictions

- We observe random variable  $X$ 
  - Want to make prediction about  $Y$
  - E.g.,  $X$  = stock price at 9am,  $Y$  = stock price at 10am
  - Let  $g(X)$  be function we use to predict  $Y$ , i.e.:  $\hat{Y} = g(X)$
  - Choose  $g(X)$  to minimize  $E[(Y - g(X))^2]$
  - Best predictor:  $g(X) = E[Y | X]$
  - Intuitively:  $E[(Y - c)^2]$  is minimized when  $c = E[Y]$ 
    - Now, you observe  $X$ , and  $Y$  depends on  $X$ , then use  $c = E[Y | X]$
- You just got your first baby steps into Machine Learning
  - We'll go into this more rigorously in a few weeks

## Speaking of Babies...

- Say my height is  $X$  inches ( $x = 71$ )
  - My son:  He does not look like: 
  - Say, historically, sons grow to heights  $Y \sim N(X + 1, 4)$ , where  $X$  is height of father
    - $Y = (X + 1) + C$  where  $C \sim N(0, 4)$
  - What should I predict for the eventual height of my son?
  - $E[Y | X = 71] = E[X + 1 + C | X = 71]$   
 $= E[72 + C] = E[72] + E[C] = 72 + 0$   
 $= 72$  inches