## Data Science Example - social ratings

Goal: Learn how statistics and probability can inform a data-driven service

Start w/ projector:

1. Reminder of expectation (mean) and variance, properties [if omitted from previous lecture]

2. Motivating slide - Amazon.com star ratings

3. Switch to board

Question: How to <u>rank</u> products using ratings of users when some products have fewer ratings than others.

Ex        Product A                          Product B
          5 out of 5 stars ($R_A = 5$)        4.5 out of 5 stars ($R_B = 4.5$)
          1 rating ($n_A = 1$)                 30 ratings ($n_B = 30$)

→ naive to sort based on R, b/c A has so few ratings.

Brainstorm solutions:      Gen Q: How to sort using two objectives

Here's some of my ideas:

• Ignore it, sort using R → bad idea

• Define <u>cutoff</u>, only show products w/ $n > n_{cutoff}$.
  → bad idea: biased to popular products, how do new products get any traction?

• Show shoppers a fancy 2D visualization    $n_i$ 
  → bad idea: way too complicated!                  $r_i$

• Scalarize the objectives to sort using one number
  → data-driven scalarization?
  → design scalarization using statistical inference ✳

# Outline

1. Problem formulation

2. Modeling a user rating

3. Modeling a product's rating

4. Connecting models to sorting problem (next time)

# 1. Problem formulation

Product $i$ has $n_i$ ratings

Design choice: 1-5 stars is complicated, <u>simplify to</u> <u>thumbs up/down</u>

- rating is now proportion of thumbs ups (t.u.)

- Real ratings often support this ⟹

 easy to transform

→ 10 people rate product, 6 give t.u. ⟹ $R_i = \frac{6}{n_i} = \frac{6}{10} = 0.6$

- This simplification is plausible for many data and makes calculations much easier, but going from a 1-5 scale to a 0-1 scale does not solve our problem of how to sort products.

Q: $R_i = \# t.u/n_i$ is the <u>observed</u> (sample) rating. What will the rating be if <u>everyone</u> (population) rated product $i$?

→ Asking this is <u>critical</u>. If we somehow knew the population rating we would know the <u>true</u> way to sort products!

But we don't know pop. rating. What information do we know about pop. rating given sample?

## 2. Modeling a user rating.

The process by which a user picks thumbs up vs. down is <u>complicated</u> and will depend on many details specific to the user and unknown to us.

A statistical model replaces these unknown factors w/ a reasonable approximation using randomness

Introduce a random variable (R.V.):

$$X_j = \begin{cases} 1 & \text{if } j^{\underline{th}} \text{ person gives t.u.} \\ 0 & \text{otherwise} \end{cases}$$

(Suppressing index of product b/c we are only considering a single product here.)

RV needs an associated prob. distribution:

$$Pr(X_j = 1) = p, \qquad Pr(X_j = 0) = 1-p$$

$$\text{or} \quad Pr(X_j = x_j) = \begin{cases} p & \text{if } x_j = 1 \\ 1-p & \text{if } x_j = 0 \end{cases}$$

• This is called a <u>Bernoulli</u> R.V. Essentially, a "coin flip"

• We have introduced a <u>parameter</u> $p$, if $p = \frac{1}{2}$ the user is equally likely to vote t.u. or t.d.

## Statistics of Bernoulli:

$$E[X_j] = \sum_x x \cdot Pr(X=x) = 1 \cdot Pr(X=1) + 0 \cdot Pr(X=0)$$
$$= 1 \cdot p + 0(1-p)$$
$$= p$$

$$Var(X_j) = \underbrace{E[X_j^2]}_{=p} - \underbrace{E[X_j]^2}_{} = \sum_x x^2 Pr(X=x) - p^2$$
$$= 1^2 \cdot p + 0^2(1-p) - p^2 = p - p^2$$
$$= p(1-p)$$

<u>Sample mean</u>: $\overline{X} = \frac{1}{n}\sum_{j=1}^n X_j$    <u>Sample variance</u>: $S_x^2 = \frac{1}{n}\sum_{j=1}^n (X_j - \overline{x})^2$

= fraction of 1's in sample

↑ be careful $n$ vs. $n-1$   [3]

# 3. Modeling a product's rating

A product will receive $n$ ratings (supressing index $i$), so we need to deal w/ **multiple** Bernoulli R.V.'s

Let's make some simplifying assumptions (that may or may not be plausible)

1. **Each** user $j$ rating the product follows the same prob. distribution

2. The parameter $p$ is **constant** for the **product** (of course, different products can and will have different $p$'s)

3. Each user rates the product **independent** of any other user's rating.

   $\rightarrow$ may not be reasonable in practice. Network effects, for example, may induce a dependency ("Everyone else loves this, so it must be great!")

Taken together, these assumptions tell us that the <u>set of $n$ ratings given to the product</u> are <u>independent and identically distributed</u> (iid). Denoted

   · The $X_j$'s are iid     · · The $\{X_j\}$ are iid.

A product's rating is the average of the $n$ i.u./i.i.d. user ratings it receives.

$$rating = \bar{X} = \frac{1}{n} \sum_{j=1}^{n} X_j$$

This is proportional to a sum:

$$n\bar{X} = \sum_{j=1}^{n} X_j \equiv k \quad \text{(call the sum } k\text{)}$$

Q we need to ask: What is the probability that the sum of $n$ "coin flips" is $k$?

A: under our assumptions, summing Bernoulli variables
gives a __binomial__ random variable

$$Pr(k; n, p) = \binom{n}{k} p^k (1-p)^{n-k}$$

R.V.  semi   parameters
     colon

idea: k 1's occur w/ prob $p^k$
      n-k 0's occur w/ prob $(1-p)^{n-k}$
      $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ arrangements of the k 1's and n-k 0's.

Understanding k lets us see how much the __collective__
rating X can vary if/when different sets of
users rate the product

## Statistics of Binomial R.V.

• $E[k] = \sum_k k \, Pr(k; n, p) = [\text{lot of work}] = np$

—or—  $E[k] = E[\sum_{j=1}^{n} X_j] = \sum_{j=1}^{n} E[X_j] = \sum_{j=1}^{n} p = np$

b/c E[ ]
is linear.

• $Var(k) = Var(\sum_{j=1}^{n} X_j) = \sum_{j=1}^{n} Var(X_j)$   (b/c $X_j$ are iid → indep)

$$= \sum_{j=1}^{n} p(1-p) = np(1-p)$$

But we really want statistics of $\bar{X}$ not k.   $(\bar{x} = \frac{k}{n})$

• $E[\bar{X}] = E[\frac{k}{n}] = \frac{1}{n} E[k] = \frac{1}{n} np = p$   makes sense!

b/c n
const.

nP

• $Var(\bar{X}) = Var(\frac{k}{n}) = \frac{1}{n^2} Var(k) = \frac{1}{n^2} np(1-p) = \frac{1}{n} p(1-p)$

careful                                          smaller✱

→ $Var(\bar{X}) < Var(X_j)$ and $Var(\bar{x})$ decreases w/ n !  makes sense: averaging
over more data (higher
n) gives less fluctuating

⇒ Recall standard __deviation__ vs. standard __error__ (or S.E.M.

⑤