Modeling Task $f = f(I, R)$

* **Statistical Models vs. Mathematical Models**
  **(Machine Learny)**

- Statistical/Machine Learny Models are fit from historical data

- Mathematical Models are equations written on paper

## Models fit from historical data

- What data do we need?
  Can we get it?
  get it.

- Choose the model

## Mathematical Models

- What Random variables/distribution

# are appropriate, if any?

- Statistical/Machine learning Models are fit from historical data

Models fit from historical data

- What data do we need? Can we get it? get it.
- Choose the model

$$f = f(I, R) = \text{Starting pit allows R runs through I innings, win prob.}$$

## data

every starting-pitcher-game $\rightsquigarrow$ index $i = 1, \ldots, n$

$R_i = $ runs allowed by s.p. in game $i$

$I_i = $ innings pitched by s.p. in game $i$

(for simplicity assume s.p. pulled after fully completing $I_i$)

$y_i = 1$ if s.p. team wins, else $0$

Lahman $\rightarrow$ box score

Retrosheet/Statcast $\Rightarrow$ PBP

# empirical gRid

$$\hat{f}(I, R) = \text{Mean}\{y_i : I = I_i \text{ and } R = R_i\}$$

$$(I, R)$$

$$(4, 2) \quad \begin{array}{l} 70\% \text{ Win} \\ 30\% \text{ Loss} \end{array} \approx .7$$

## XGBoost with Monotonic Constraints

XGBoost = one of the fastest and
easiest to use "off-the-shelf"
machine learning algorithms
supervised

dependent variables $X_1, ..., X_p$
independent variables $y$

XGBoost "memorizes" as best as
possible
$$g(x) = y \longrightarrow \text{interpolates}$$

Monotone Constraints:

$$R \longrightarrow tell \; X_{(about)} \; to \; be \; \underset{in \; R}{\overset{monotonic}{decreasing}}$$

$$I \longrightarrow tell \; X_{(about)} \; to \; be \; \underset{increasing \; in \; I}{monotonic}$$

## Try Mathematical Models

$(R, I) \longrightarrow$ win probability

say inning $i$   $X_i = \#$ Runs scored by team the pitcher's

$Y_i = \#$ Runs scored by the opp's team

$$\hat{f}(I, R) = \mathbb{P}\left( \sum_{i=1}^{9} X_i > R + \sum_{i=I+1}^{9} Y_i \right)$$

$$+ \frac{1}{2} \mathbb{P}\left( \sum_{i=1}^{9} X_i = R + \sum_{i=I+1}^{9} Y_i \right)$$

$\boxed{\text{Start simple}} \longrightarrow$ Poisson
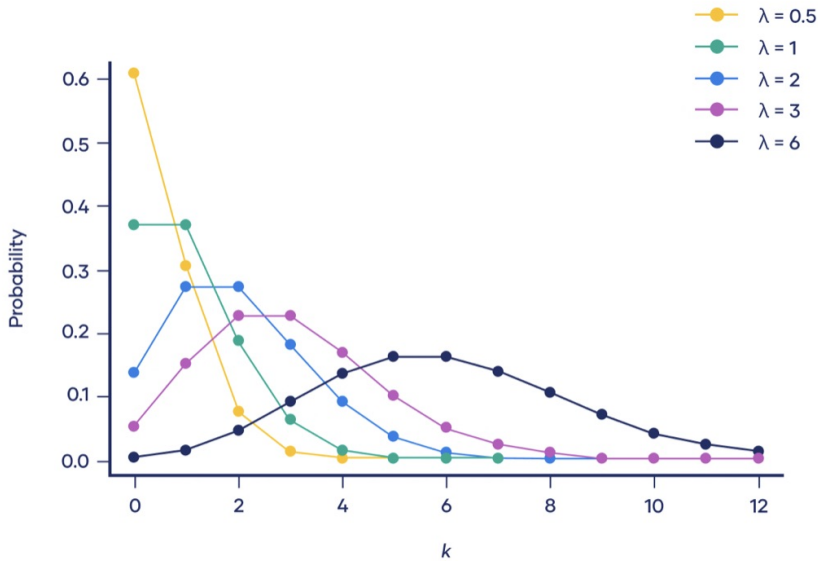
$X_i, Y_i$ random variables (distributions)

$\hookrightarrow$ possible outcomes $0, 1, 2, 3, \ldots$

$$X \sim Poisson(\lambda) \quad means$$

Outcomes $\quad x \in \{0, 1, 2, 3, \ldots\}$

$$\mathbb{P}(X = x) = e^{-\lambda x} \frac{\lambda^x}{x!}$$

$$\mathbb{E}X = \lambda, \quad var(X) = \lambda, \quad \lambda > 0$$



$$\mathbb{P}\left( \sum_{i=1}^{q} X_i > R + \sum_{i=I+1}^{q} Y_i \right)$$

$$= \left\{ \mathbb{P}\left( \sum_{i=1}^{q} Poisson(\lambda) > R + \sum_{i=I+1}^{q} Poisson(\lambda) \right) \right.$$

$$\text{if } I < q$$

$$\mathbb{P}\left( \sum_{i=1}^{q} \text{Poisson}(\lambda) > R \right) \qquad \text{if } I = 9$$

$$= \begin{cases} \mathbb{P}\left[ \text{Poisson}(q \cdot \lambda) > R + \text{Poisson}\left[(q - (I+1))\lambda\right] \right] \\ \\ \mathbb{P}\left( \text{Poisson}(q \cdot \lambda) > R \right) \qquad \text{if } I = 9 \end{cases}$$

__Thm__  $\text{Poisson}(\lambda_1) + \text{Poisson}(\lambda_2) = \text{Poisson}(\lambda_1 + \lambda_2)$

$$= \begin{cases} \mathbb{P}\left[ \text{Skellam}(q\lambda, (q - I - 1)\lambda) > R \right] \quad \text{if } I < q \\ \\ \mathbb{P}\left[ \text{Poisson}(q\lambda) > R \right] \qquad \text{if } I = q \end{cases}$$

__Thm__  $\text{Poisson}(\lambda_1) - \text{Poisson}(\lambda_2)$
$\qquad = \text{Skellam}(\lambda_1, \lambda_2)$

Now have an explicit formula for
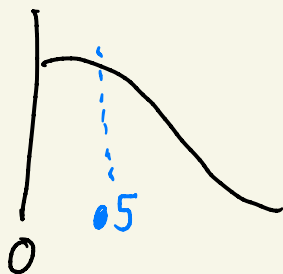$f(I,R)$ if we assume

$$X_i, Y_i \sim Poisson(\lambda).$$

$$\hat{\lambda} = \mathbb{E}X_i = \text{mean runs allowed in}$$
$$\text{an inning by one}$$
$$\text{team}$$

Need to choose a smart value of $\lambda$

for a given league-season (e.g. 2019 NL)
let $\lambda =$ observed mean runs allowed
in a half-inning

Code

$$X_i \sim \text{Poisson}(\lambda_x)$$

$$Y_i \sim \text{Poisson}(\lambda_y)$$

$$\lambda_x^{(m)}, \lambda_y^{(m)} \sim \mathcal{N}_+(\lambda, \sigma^2)$$

$$f(R, I \mid \lambda_x, \lambda_y) =$$

$$\mathbb{P}\left( \sum_{i=1}^{q} X_i > R + \sum_{i=I+1}^{q} Y_i \right)$$

$$+ \frac{1}{2} \mathbb{P}\left( \sum_{i=1}^{q} X_i = R + \sum_{i=I+1}^{q} Y_i \right)$$

= a similar formula as before
except now in terms of $\lambda_x, \lambda_y$

$$f(I, R) = \frac{1}{M} \sum_{m=1}^{M} f\left(I, R \mid \lambda_x^{(m)}, \lambda_y^{(m)}\right)$$

$$\begin{cases} X_i \sim \text{Poisson}(\lambda_x) \\ Y_i \sim \text{Poisson}(\lambda_y) \\ \lambda_x^{(m)}, \lambda_y^{(m)} \sim N_+(\lambda, \sigma^2) \end{cases}$$

$$\lambda = \text{mean}\left( \begin{Bmatrix} \text{mean} \\ \text{Runs scored in} \\ \text{half inning by} \\ \text{team } t \end{Bmatrix} \right)$$

$$\sigma^2 = \text{var}\left( \left\{ \; \uparrow \; \right\} \right)$$

2019 NL: get $\hat{\lambda}, \hat{\sigma}^2$

then

$$f(I, R \mid \hat{\lambda}, \hat{\sigma}^2) = \frac{1}{M} \sum_{m=1}^{M} f\left(I, R \mid \lambda_x^{(m)}, \lambda_y^{(m)}\right)$$

where $\lambda_x^{(m)}, \lambda_y^{(m)} \sim N_+(\hat{\lambda}, \hat{\sigma}^2)$

$M = 100$

$$\begin{cases} X_i \sim \text{Poisson}(\lambda_x) \\ Y_i \sim \text{Poisson}(\lambda_y) \\ \lambda_x^{(m)}, \lambda_y^{(m)} \sim \mathcal{N}_+(\lambda, \sigma^2 \cdot K) \end{cases}$$

$$K < 1$$

$$f(I, R \mid \hat{\lambda}, \sigma^2, K) = \frac{1}{M} \sum_{m=1}^{M} f\left(I, R \mid \lambda_x^{(m)}, \lambda_y^{(m)}\right)$$

$$\text{where } \lambda_x^{(m)}, \lambda_y^{(m)} \sim \mathcal{N}_+\left(\hat{\lambda}, \sigma^2 \cdot K\right)$$

$\hookrightarrow$ grid $f(I, R \mid \hat{\lambda}, \sigma^2, K)$

how to choose $K$?

try different $K$'s and choose the $K$ which works best, meaning makes the most accurate predictions.

WP predictions $f(I, R | \hat{\lambda}, \hat{\sigma}^2, K)$

observed Win/loss column

logloss $(pred, obs. win/loss)$