Ex : Sparrows

Y: # fledglings (offspring)

X: ages

in a mating season

idea: The # of Fledglings a sparrow has ^ depends
on the age of the sparrow.

We could model this count data Y with :

$$Y|X \sim Poisson(\theta_x)$$ since support Poisson r.v. Y is $\{0,1,2,\ldots\}$

Here, $\mathbb{E}Y|X = \theta_x$ is age-specific.

Possible ways to model $\theta_x$:

(1) discrete: $\theta_1, \theta_2, \ldots, \theta_6$
   a specific $\theta_i$ for each age $i$ of sparrow.

* Note: if we don't collect much data on certain
   age sparrows then our estimates for some $\theta_i$
   will be poor.

(2) $\theta_x = f(x) = \beta_1 + \beta_2 X + \beta_3 X^2$
   _equation of a hyperplane_

   * Note $\theta_x$ should _not_ be negative yet it _can_ be
      under model (2).

(3) sol'n : log-transform.

$$\log \mathbb{E}Y|X = \log \theta_{xi}: \quad \beta_1 + \beta_2 X_i + \beta_3 X_i^2$$

i.e. $\mathbb{E}Y|X = \exp\{\beta_1 + \beta_2 X + \beta_3 X^2\} > 0$

Terminology

$$Y|X \sim Poisson\left(\exp\{\beta^T \vec{x}\}\right)$$   "Poisson regression"

$\beta^T x$                                    "linear predictor"

$\mathbb{E}Y|X$ is <u>linked</u> to the linear predictor w/ the <u>log</u> function. So we say this model has a <u>log-link</u>.

In general if $f(\mathbb{E}Y|X) = \beta^T X$ is a <u>generalized linear model</u> (GLM) w/ link function $f$.

---

A note on $\beta^T \vec{x}$

$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$    Dimension $3 \times 1$

$\vec{x}_i = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix}$ Dimension $3 \times 1$    i.e. corresponds to 1 sparrow.

So for any individual sparrow,

$$E y_i | x_i = \exp\{\underset{1\times3\ \ 3\times1}{\beta^T x_i}\} = \exp\{\beta_1 \cdot 1 + \beta_2 x_i + \beta_3 x_i^2\}$$

---

Recap:

We have a data generative model:

$$Y|X \sim \text{Poisson}(\exp\{\beta^T x\})$$

and if we know a bird's age <u>and</u> <u>we knew $\beta$</u>, we <u>could predict</u> $Y$.

The trouble is, <u>we don't know $\beta$</u>. $\beta_1, \beta_2, \beta_3$ are unknown. We want to estimate them from some data. We want $P(\beta | Y, X)$. We need to specify priors on $\beta$.

Common prior:

$$\beta \sim MVN(0, \Sigma)$$

EX    What does $p(\beta \mid y, x)$ look like?

Sol

$$p(\beta \mid y, x) \propto p(\beta) \prod_{i=1}^{n} p(y_i \mid \beta, x_i)$$

$$\propto \underbrace{\exp\left\{-\tfrac{1}{2}\beta^T \Sigma \beta\right\}}_{\text{normal}} \underbrace{\prod_{i=1}^{n} \exp\left\{(\beta^T x_i) \cdot y_i\right\} \cdot \exp\left\{-\exp\left\{\beta^T x_i\right\}\right\}}_{\exp\left\{\sum_{i=1}^{n}\left[\beta^T x_i y_i - \exp\{\beta^T x_i\}\right]\right\}}$$

$$\overset{?}{\text{Does not look like a known kernel.}}$$

Rec    NOT    conjugate

NOT    semi-conjugate

NOT    easy to sample from using known methods.

But still, we want to sample from it because

$$p(\beta \mid y, x) = \frac{p(\beta) \, p(y \mid x, \beta)}{\iiint p(\beta) \, p(y \mid x, \beta) \, d\beta_1 \, d\beta_2 \, d\beta_3} \overset{\rightarrow}{\leftarrow} "p(y)"$$

is a nasty integral to numerically approx in higher dimensions.

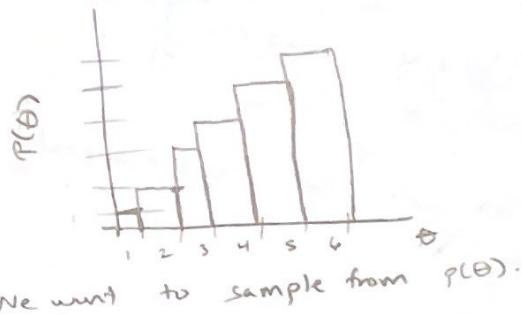Goal: construct a Markov chain that approximates the posterior. we need to __avoid__ computing $p(y)$.

Metropolis algo.

<u>Intuition</u> :  weighted die ex.

Let $p(\theta = i) = i/w$   for   $i \in \{1, ..., 6\}$

So the pmf of $\theta$ looks like:



We want to sample from $p(\theta)$.

Start @ $\theta = 3$

propose $\theta^* = 2$

We need more $\theta = 3$ than $\theta = 2$ in our sample, if our sample is to approximate $p(\theta)$ well. How many more?

$$\frac{p(\theta^*)}{p(\theta)} = \frac{2/w}{3/w} = \frac{2}{3} \cdot$$   I need $2/3$ as many samples of $\theta = 2$ in my Markov chain.

So accept new state $\theta^* = 2$ w/ probability $2/3$.

Say, I propose $\theta^* = 4$. Current state is $\theta = 3$

$$\frac{p(\theta^*)}{p(\theta)} = \frac{4/w}{3/w} > 1 \cdot$$   So the state $\theta = 4$ is more probable. We should accept it as the next state in our Markov chain.

Eventually we will have:

$$\theta \begin{bmatrix} 3 \\ 2 \\ 3 \\ 4 \\ 4 \\ 5 \\ 4 \\ 5 \\ 6 \\ 5 \\ \vdots \end{bmatrix}$$

that will approximate the target distr. well.

* <u>Note</u> we do <u>not</u> need to know $w$ to implement this.

## Metropolis Algorithm

1. Sample $\theta^* | \theta^{(s)} \sim J(\theta | \theta^{(s)})$

2. Compute acceptance ratio

$$r = \frac{p(\theta^* | y)}{p(\theta | y)}$$

3. Let $\theta^{(s+1)} = \begin{cases} \theta^* & w/ \quad prob \quad \min(r, 1) \\ \theta^{(s)} & w/ \quad prob \quad 1 - \min(r, 1) \end{cases}$

$J(\theta | \theta^{(s)})$ is the _proposal_ distribution. It proposes a new value $\theta$ _given_ our current $\theta^{(s)}$.

For this to be the "Metropolis algo.", $J$ is _symmetric_.
i.e. $J(\theta_a | \theta_b) = J(\theta_b | \theta_a)$.

_Practice_     (i) implement die example, w/
$J(\theta = j | \theta^{(s)} = i) = 1/6$ for all $j$ :
i.e. propose a new state $j$ uniformly.