

# Count data and the Poisson distribution

---

Frank Edwards

- Counts are cumulative totals of the number of incidences of some event, generally across time or place

- Counts are cumulative totals of the number of incidences of some event, generally across time or place
- Counts are positive integers  $\in [0, \infty]$

- Counts are cumulative totals of the number of incidences of some event, generally across time or place
- Counts are positive integers  $\in [0, \infty]$

## Counts as extensions of binary data

- Counts can be thought of as repeated binary trials
- $\sum y_i$  where  $y$  is equal to 1 or 0 provides a count
- Generally, we could treat `sum(y==1) + sum(y==0)` or `nrow(y)` as the exposure, or denominator for a rate. Why?

# An example of count data

```
load("../data/fieldplayer_overall_season_stats.rda")
load("../data/player.rda")
```

```
nwsl_stats<-fieldplayer_overall_season_stats
nwsl_players<-player
```

```
head(nwsl_players)
```

```
## # A tibble: 6 x 5
##   person_id player      nation pos  name_other
##   <dbl> <chr>      <chr> <chr> <chr>
## 1     342 Marisa Abegg    USA   DF   <NA>
## 2     117 Danesha Adams  USA   FW,MF <NA>
## 3        6 Adriana      ESP   FW   <NA>
## 4     300 Leigh Ann Brown USA   DF,MF <NA>
## 5     202 Jazmyne Avant  USA   DF   <NA>
## 6      28 Amy Barczuk   USA   DF   <NA>
```

## Approaches to modeling count data

---

# The Poisson model

Where  $y$  is a non-negative integer (count)

$$y \sim \text{Poisson}(\lambda)$$

$$E(y) = \bar{y} = \lambda$$

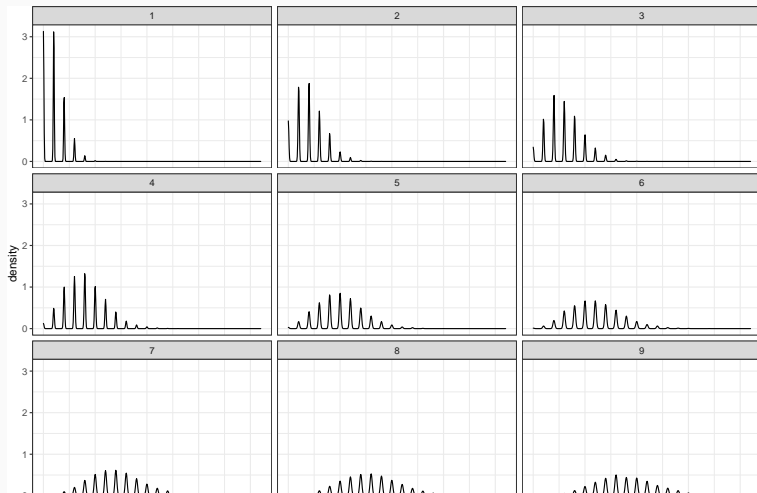
$$\text{Var}(y) = \lambda$$

$$\text{Pr}(y = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$



# Shape of the Poisson distribution

```
ggplot(pois_demo, aes(x=count)) +  
  geom_density(adjust = 1/4) +  
  facet_wrap(~lambda)
```



## Let's look at each Poisson variable

```
pois_demo%>%group_by(lambda)%>%  
  
  summarise(mean = mean(count),  
            variance = var(count))
```

```
## # A tibble: 9 x 3  
##   lambda mean variance  
## *   <int> <dbl>     <dbl>  
## 1     1  1.01     1.01  
## 2     2  2.00     1.99  
## 3     3  2.97     2.94  
## 4     4  4.00     3.99  
## 5     5  5.02     5.04  
## 6     6  6.02     6.27  
## 7     7  7.00     6.93  
## 8     8  7.98     8.11  
## 9     9  8.98     9.20
```

For a count variable  $y$ , we can specify a Poisson GLM with a log link function

$$y \sim \text{Poisson}(\lambda)$$

$$\lambda = \beta X = \beta_0 + \beta_1 x_1 \cdots \beta_n x_n$$

$$E(y|x) = e^\lambda$$

$$\log(E(y|x)) = \lambda = \beta X$$

## Advantages of the Poisson distribution for regression

1. Constrained to non-negative integers
2. Variance scales with the expectation of  $y$
3. Relatively simple to interpret

However:

$$\lambda = E(y|x) = \text{var}(y)$$