# Generalized Linear Models

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$y = X^T \beta + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$

$$\approx \quad y \sim N(X^T \beta, \sigma^2)$$

$X$ wing span

$Y$

Lim beak length $\sim N(\mu, \beta^2)$

GLM Population size $\sim Pois(\lambda)$

Like talon size $\sim N(\mu, \sigma^2)$

Species

Can predictors

explain the outcome?

Logistic Reg

GLM

$$E(Y) = g^{-1}(X^T\beta)$$

10cm humming bird $\sim$ Bern($p$)

20cm Seagull

30cm duck

1m Conor

$\sim$ Bin($n=1, p$)

$\sim$ Multinomial($p_k$)

$X = 0.85m$

$X = 25cm$

$X = 0.25m$

likelihood

$$\ell(\beta) \cdot + \lambda \beta^2$$

Prediction     vs     Explanation

Bayes Classifiers

Naive Bayes

Linear Discrim Analysis

k Nearest Neighbors

# Learning Outcomes

- Exponential Family of Distributions

- Generalized Linear Models

$$g?$$

# Exponential Family of Distributions

$$E(Y) = g^{-1}(X^T \beta)$$

# Exponential Family of Distributions

An exponential family of distributions are random variables that allow their probability density function to have the following form:

$$f(y; \theta, \phi) = a(y, \phi) \exp\left\{ \frac{y\theta - \kappa(\theta)}{\phi} \right\}$$

$a(y, \phi) =$ normalizing constant

$\phi =$ dispersion parameter function

$\longrightarrow \theta =$ canonical link parameter/function

$\kappa(\theta) =$ log-cumulant function

# Canonical Parameter

The canonical parameter represents the relationship between the random variable and the $E(Y) = \mu$

$$\eta = X^T \beta = g(\mu)$$

# Normal Distribution

$$f_{(x)} = a(x; \phi) \, e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$E(X) = \mu$$

$$a(x; \phi) \propto \frac{1}{\sqrt{2\pi\sigma^2}} \qquad \phi = \sigma^2$$

$$f_{(x)} = a(x; \phi) \, e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{1}{2\sigma^2}\left[x^2 - 2x\mu + \mu^2\right]}$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{x^2}{2\sigma^2} + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}}$$

$$f = a(\cdot) \, e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$a(x; \phi) = \frac{e^{-\frac{x^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} \cdot e^{\frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}}$$

$$\frac{e^{-x^2/2\sigma^2}}{\sqrt{2\pi\sigma^2}} \qquad e^{\frac{x\mu - \mu^2/2}{\sigma^2}}$$

$$f = a(\cdot) e^{\frac{x\theta - h(\theta)}{\phi}}$$

$$E(Y) = \beta_0 + \beta_1 x \cdots$$

$$\theta = \mu = X^T\beta \qquad E(Y) = X^T\beta$$

$$\cdot E(Y) = X^T\beta$$

# Binomial Distribution

$$f_{(x)} = a(x;\phi) e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$f(x; n, p) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$E(y) = p$$

$$a(x;\phi) \propto \binom{n}{x}$$

$$y = e^{\ln y}$$

$$\binom{n}{x} e^{\ln\left(p^x (1-p)^{n-x}\right)}$$

$$\binom{n}{x} e^{\ln(p^x) + \ln\left[(1-p)^{n-x}\right]}$$

$$f = a(\cdot) e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$\binom{n}{x} e^{x\ln p + (n-x)\ln(1-p)}$$

$$\binom{n}{x} e^{x(\ln p - \ln(1-p)) + n\ln(1-p)}$$

$$\binom{n}{x} e^{\frac{x\ln\left(\frac{p}{1-p}\right) + n\ln(1-p)}{1}}$$

$$\binom{n}{x} e^{\frac{x\ln\left(\frac{p}{1-p}\right) + n\ln(1-p)}{1}}$$

$$f = a(\cdot) e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$\phi = 1$$

logit function

$$\implies \ln\left(\frac{p}{1-p}\right) = X^T\beta$$

$$E(Y) = p = \frac{e^{X^T\beta}}{1 + e^{X^T\beta}}$$

$$E(Y)$$

$$\frac{(P=1)}{(P=0)} \quad \begin{matrix} \text{odds} \\ \downarrow \\ \text{success} \\ \text{failure} \end{matrix}$$

Logistic Regression

$$Y = \begin{matrix} 1 \\ 0 \end{matrix}$$

$$X: X_{,,} \ldots, X_p$$

# Poisson Distribution

$$f(\lambda) = a(x; \phi) \, e^{\frac{x\theta - k(\theta)}{\phi}}$$

$$f(x; \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\widehat{E(Y)} = \lambda$$

$$a \propto \frac{1}{x!} e^{-\lambda} e^{\ln(\lambda^x)}$$

$$\frac{1}{x!} e^{-\lambda + \ln(\lambda^x)} = \frac{1}{x!} e^{\frac{x\ln(\lambda) - \lambda}{1}}$$

$$\theta = \ln(\lambda) = X^T \beta$$

# Common Distributions and Canonical Parameters

| Random Variable | Canonical Parameter $\theta$ | |
|---|---|---|
| Normal | $\mu$ | Identity link |
| Binomial | $\log\left(\frac{\mu}{1-\mu}\right)$ | logit |
| Negative Binomial | $\log\left(\frac{\mu}{\mu+k}\right)$ | logit |
| Poisson | $\log(\mu)$ | log ln |
| Gamma | $-\frac{1}{\mu}$ | inverse $\frac{1}{\mu}$ |
| Inverse Gaussian | $-\frac{1}{2\mu^2}$ | ~ double inverse |

| Random Variable | Canonical Parameter |
| --- | --- |

# Generalized Linear Models

# Generalized Linear Models

A generalized linear model (GLM) is used to model the association between an outcome variable (of any data type) and a set of predictor values. We estimate a set of regression coefficients $\beta$ to explain how each predictor is related to the expected value of the outcome.

$$E(Y) = g^{-1}(X^T\beta)$$

$$X^T\beta = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

# Generalized Linear Models

A GLM is composed of a systematic and random component.

$$Sys: \quad X^T \beta$$

$$Random: \quad Distribution$$

# Random Component

The random component is the random variable that defines the randomness and variation of the outcome variable.

# Systematic Component

The systematic component is the linear model that models the association between a set of predictors and the expected value of Y:

$$g(\mu) = \eta = X_i^{\mathrm{T}}\boldsymbol{\beta}$$

Function RV    $X \sim$ Normal

$$Y = X^2$$         what is $Y$?

$X \sim$ Pois

Var$(Y)$    $Y = 2X + 5$

Central Limit Theorem    $n \rightarrow \infty$

Sampling Distribution    $\bar{X}$ ? $S^2$ ?

$X \sim$ Normal

MLE    $\Leftarrow$

Linear Regression $\Leftarrow$ do math

GLM's     Link Functions,

Estimation ⟵