

Learning From Data

Lecture 3: Generalized Linear Models

Yang Li yangli@sz.tsinghua.edu.cn

September 30, 2022

Ask me a question

What is the difference between probabilistic and non-probabilistic methods?

probabilistic: $h_{\theta}(x) = \begin{cases} 0.9 & \text{if } \frac{1}{2} \sigma(\sigma^T x) \geq \Pr(y=1|x) \\ 0.1 & \text{if } \Pr(y=0|x) \end{cases} \leftarrow \text{advantage?}$

clf. predict-probal

if $\Pr(y=1|x) > \Pr(y=0|x) \Rightarrow$ label $y=1$

$y = \begin{cases} 1 & \Pr(y=1|x) > 0.5 \\ 0 & \text{o.w.} \end{cases}$

Non-probabilistic


(deterministic) $h_{\theta}(x) = \text{sign}(w^T x + b) = \begin{cases} 1 & w^T x + b > 0 \\ 0 & \text{o.w.} \end{cases}$

clf. predict()

- gives the confidence of the prediction
- better understanding of the model and data

Today's Lecture

Supervised Learning (Part III)

- ▶ Review on linear and logistic regression
 - ▶ Softmax Regression
 - ▶ Review: exponential families
 - ▶ Generalized linear models (GLM)
- 

Written Assignment (WA1) is released. Due on Oct 8th. (Start early!)

Review of Lecture 2

Review of Lecture 2: Linear least square

- ▶ Hypothesis function for input feature $x^{(i)} \in \mathbb{R}^n$:

$$\underline{h_\theta(x^{(i)})} = \underline{\theta^T x^{(i)}}, \text{ where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

Review of Lecture 2: Linear least square

- ▶ Hypothesis function for input feature $x^{(i)} \in \mathbb{R}^n$:

$$h_{\theta}(x^{(i)}) = \theta^T x^{(i)}, \text{ where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

- ▶ Cost function for m training examples $(x^{(i)}, y^{(i)}), i = 1, \dots, m$:

$$J(\theta) =$$

Review of Lecture 2: Linear least square

- ▶ Hypothesis function for input feature $x^{(i)} \in \mathbb{R}^n$:

$$h_{\theta}(x^{(i)}) = \theta^T x^{(i)}, \text{ where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

- ▶ Cost function for m training examples $(x^{(i)}, y^{(i)}), i = 1, \dots, m$:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

$\frac{1}{n} \sum (\theta)$

Also known as **ordinary least square regression** model.

Review of Lecture 2: Linear least square

- ▶ Hypothesis function for input feature $x^{(i)} \in \mathbb{R}^n$:

$$h_{\theta}(x^{(i)}) = \theta^T x^{(i)}, \text{ where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

- ▶ Cost function for m training examples $(x^{(i)}, y^{(i)}), i = 1, \dots, m$:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left(y^{(i)} - \theta^T x^{(i)} \right)^2$$

Also known as **ordinary least square regression** model.

How to minimize $J(\theta)$?

▶ Gradient descent:

update rule (batch)

update rule (stochastic)

▶ Newton's method

▶ Normal equation

How to minimize $J(\theta)$?

- ▶ Gradient descent:

update rule (batch) $\theta_j \leftarrow \theta_j + \alpha \cdot \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$

Handwritten notes: $-\nabla J(\theta_j)$ above the equation, and a green box around the entire update rule.

update rule (stochastic)

- ▶ Newton's method

- ▶ Normal equation

How to minimize $J(\theta)$?

- ▶ Gradient descent:

$$\text{update rule (batch)} \quad \theta_j \leftarrow \theta_j + \alpha \cdot \frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)}$$

$$\text{update rule (stochastic)} \quad \theta_j \leftarrow \theta_j + \alpha \left(y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)}$$

- ▶ Newton's method

- ▶ Normal equation

How to minimize $J(\theta)$?

- ▶ Gradient descent:

$$\text{update rule (batch)} \quad \theta_j \leftarrow \theta_j + \alpha \cdot \frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)}$$

$$\text{update rule (stochastic)} \quad \theta_j \leftarrow \theta_j + \alpha \left(y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)}$$

- ▶ Newton's method

$$\nabla J(\theta) = 0.$$

$$\theta \leftarrow \theta - H^{-1} \nabla J(\theta)$$

- ▶ Normal equation

$$\underline{X^T X \theta = X^T y}.$$

$$\theta = (X^T X)^{-1} X^T y$$

Review of Lecture 2

Maximum likelihood estimation

- ▶ Log-likelihood function:

$$\ell(\theta) = \log \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) \right) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

where p is a probability density function.

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

Review of Lecture 2

Maximum likelihood estimation

- ▶ Log-likelihood function:

$$\ell(\theta) = \log \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) \right) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

where p is a probability density function.

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

(True or False?) Ordinary least square regression is equivalent to the maximum likelihood estimation of θ .

Review of Lecture 2

Maximum likelihood estimation

- ▶ Log-likelihood function:

$$\ell(\theta) = \log \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) \right) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

where p is a probability density function.

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

(True or False?) Ordinary least square regression is equivalent to the maximum likelihood estimation of θ .

True under the assumptions:

- ▶ $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$
- ▶ $\epsilon^{(i)}$ are i.i.d. according to $\mathcal{N}(0, \sigma^2)$

Review of Lecture 2: Logistic regression

- ▶ Hypothesis function:

$$h_{\theta}(x) = g(\theta^T x), \quad g(z) = \frac{1}{1 + e^{-z}} \text{ is the sigmoid function.}$$

Review of Lecture 2: Logistic regression

- ▶ Hypothesis function:

$$h_{\theta}(x) = g(\theta^T x), \quad g(z) = \frac{1}{1 + e^{-z}} \text{ is the sigmoid function.}$$

- ▶ Assuming $\underline{y|x; \theta}$ is distributed according to Bernoulli($h_{\theta}(x)$)

$$p(\underline{y|x; \theta}) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$$

Review of Lecture 2: Logistic regression

- ▶ Hypothesis function:

$$h_{\theta}(x) = g(\theta^T x), \quad g(z) = \frac{1}{1 + e^{-z}} \text{ is the sigmoid function.}$$

- ▶ Assuming $y|x; \theta$ is distributed according to $\text{Bernoulli}(h_{\theta}(x))$

$$p(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$$

Review of Lecture 2: Logistic regression

- ▶ Hypothesis function:

$$h_{\theta}(x) = g(\theta^T x), \quad g(z) = \frac{1}{1 + e^{-z}} \text{ is the sigmoid function.}$$

- ▶ Assuming $y|x; \theta$ is distributed according to Bernoulli($h_{\theta}(x)$)

$$p(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$$

- ▶ Log-likelihood function for m training examples:

$$\ell(\theta) = \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))$$

Review of Lecture 2: Multi-Class Classification

Approach 1: Turn multi-class classification to a binary classification problem.

One-Vs-Rest

Learn k classifiers h_1, \dots, h_k . Each h_i classify one class against the rest of the classes.

Given a new data sample x , its predicted label \hat{y} :

$$\hat{y} = \underset{i}{\operatorname{argmax}} h_i(x)$$

Review of Lecture 2: Multi-Class Classification

Approach 1: Turn multi-class classification to a binary classification problem.

One-Vs-Rest

Learn k classifiers h_1, \dots, h_k . Each h_i classify one class against the rest of the classes.

Given a new data sample x , its predicted label \hat{y} :

$$\hat{y} = \underset{i}{\operatorname{argmax}} h_i(x)$$

Drawbacks of One-Vs-Rest:

- ▶ Class imbalance: more negative samples than positive samples when k is large

Review of Lecture 2: Multi-Class Classification

Approach 1: Turn multi-class classification to a binary classification problem.

One-Vs-Rest

Learn k classifiers h_1, \dots, h_k . Each h_i classify one class against the rest of the classes.

Given a new data sample x , its predicted label \hat{y} :

$$\hat{y} = \underset{i}{\operatorname{argmax}} h_i(x)$$

Drawbacks of One-Vs-Rest:

- ▶ Class imbalance: more negative samples than positive samples when k is large

Approach 2: Multinomial classifier (one model for all classes)

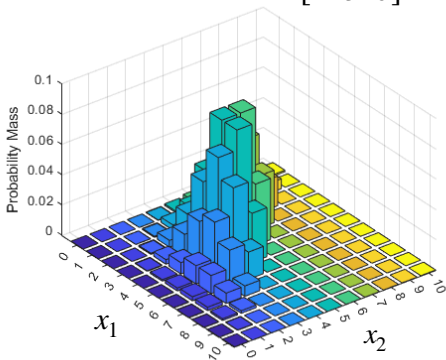
Softmax Regression

Review: Multinomial Distribution

Models the probability of counts for each side of a k -sided die rolled m times, each side with independent probability ϕ_i

$$\phi_1 + \dots + \phi_k = 1$$

$$k = 3, n = 10 \quad \phi = \left[\frac{1}{2}, \frac{1}{3}, \frac{1}{6} \right]$$



Extend logistic regression: Softmax Regression

Assume $p(y|x)$ is **multinomial distributed**, $k = |\mathcal{Y}|$

Extend logistic regression: Softmax Regression

$$\theta = [\theta_1, \dots, \theta_k]$$

$$\theta_i \in \mathbb{R}^n$$

Assume $p(y|x)$ is **multinomial distributed**, $k = |\mathcal{Y}|$

Hypothesis function for sample x :

$$\underline{h_\theta(x)} = \begin{bmatrix} p(y=1|x;\theta) \\ \vdots \\ p(y=k|x;\theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x_j}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_k^T x} \end{bmatrix} = \text{softmax}(\theta^T x)$$

$\rightarrow h_\theta(x)_i$
 $h_\theta(x)_k$

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$$p(y=i|x;\theta) = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}}$$

Extend logistic regression: Softmax Regression

Assume $p(y|x)$ is **multinomial distributed**, $k = |\mathcal{Y}|$

Hypothesis function for sample x :

$$h_{\theta}(x) = \begin{bmatrix} p(y = 1|x; \theta) \\ \vdots \\ p(y = k|x; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x_j}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_k^T x} \end{bmatrix} = \text{softmax}(\theta^T x)$$

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Parameters: $\theta = \begin{bmatrix} - & \theta_1^T & - \\ & \vdots & \\ - & \theta_k^T & - \end{bmatrix}$ } k .

Softmax Regression

Given $(x^{(i)}, y^{(i)})$, $i = 1, \dots, m$, the log-likelihood of the Softmax model is

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k p(y^{(i)} = l | x^{(i)})^{\mathbf{1}\{y^{(i)}=l\}}\end{aligned}$$

Softmax Regression

Given $(x^{(i)}, y^{(i)})$, $i = 1, \dots, m$, the log-likelihood of the Softmax model is

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k p(y^{(i)} = l | x^{(i)}) \mathbf{1}_{\{y^{(i)}=l\}} \\ &= \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}_{\{y^{(i)} = l\}} \log p(y^{(i)} = l | x^{(i)})\end{aligned}$$

Softmax Regression

Given $(x^{(i)}, y^{(i)})$, $i = 1, \dots, m$, the log-likelihood of the Softmax model is

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k p(y^{(i)} = l | x^{(i)}) \mathbf{1}\{y^{(i)}=l\} \\ &= \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}\{y^{(i)} = l\} \log p(y^{(i)} = l | x^{(i)}) \\ &= \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}\{y^{(i)} = l\} \log \frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}}\end{aligned}$$

Softmax Regression

Derive the stochastic gradient descent update: $\sum_{j=1}^K e^{\theta_j^T x}$.

- Find $\nabla_{\theta_l} \ell(\theta)$

$$\nabla_{\theta_l} \ell(\theta) = \sum_{i=1}^m \left[\left(\mathbf{1}\{y^{(i)} = l\} - P(y^{(i)} = l | x^{(i)}; \theta) \right) x^{(i)} \right]$$

Property of Softmax Regression

- Parameters $\theta_1, \dots, \theta_k$ are not independent:

$$\sum_j p(y = j|x) = \sum_j \phi_j = 1$$
- Knowing $k - 1$ parameters completely determines model.

Invariant to parameter shift

$$p(y|x; \theta) = p(y|x; \theta - \psi) \quad \theta, \psi \in \mathbb{R}^n$$

↖ $\begin{bmatrix} \theta_1 - \psi \\ \theta_2 - \psi \\ \vdots \\ \theta_k - \psi \end{bmatrix}$

Proof. $p(y=l|x; \theta - \psi)$

$$= \frac{e^{(\theta_l - \psi)^T x}}{\sum_{j=1}^k e^{(\theta_j - \psi)^T x}} = \frac{e^{\theta_l^T x} \cdot e^{-\psi^T x}}{\sum_{j=1}^k e^{\theta_j^T x} \cdot e^{-\psi^T x}} = \frac{e^{\theta_l^T x} \cdot \cancel{e^{-\psi^T x}}}{\left(\sum_{j=1}^k e^{\theta_j^T x} \right) \cdot \cancel{e^{-\psi^T x}}} = p(y=l|x; \theta)$$

Relationship with Logistic Regression

When $K = 2$,

$$\underline{h_{\theta}(x)} = \frac{1}{\underline{e^{\theta_1^T x}} + \underline{e^{\theta_2^T x}}} \begin{bmatrix} e^{\theta_1^T x} \\ e^{\theta_2^T x} \end{bmatrix}$$

Relationship with Logistic Regression

When $K = 2$,

$$h_{\theta}(x) = \frac{1}{e^{\theta_1^T x} + e^{\theta_2^T x}} \begin{bmatrix} e^{\theta_1^T x} \\ e^{\theta_2^T x} \end{bmatrix}$$

Replace $\underline{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$ with $\underline{\theta}_* = \theta - \begin{bmatrix} \theta_2 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} \theta_1 - \theta_2 \\ 0 \end{bmatrix}$,

$$\underline{h_{\theta}(x)} = \frac{1}{e^{\theta_1^T x - \theta_2^T x} + e^{0^T x}} \begin{bmatrix} e^{(\theta_1 - \theta_2)^T x} \\ e^{0^T x} \end{bmatrix}$$

$$= \left[\frac{e^{(\theta_1 - \theta_2)^T x}}{1 + e^{(\theta_1 - \theta_2)^T x}} \right]$$

$$\frac{1}{1 + e^{-z}}, z = (\theta_1 - \theta_2)^T x = \theta_*^T x$$

$$= \left[\frac{1}{1 + e^{-(\theta_1 - \theta_2)^T x}} \right] = \left[\frac{g(\theta_*^T x)}{1 - g(\theta_*^T x)} \right]$$

When to use Softmax?

$y = 1$ mammal
 $y = 2$ dog. } not mutually exclusive!

- ▶ When classes are mutually exclusive: use Softmax
- ▶ Not mutually exclusive (a.k.a. multi-label classification): multiple binary classifiers may be better

Summary: Linear models

What we've learned so far:

Learning task	Model	$p(y x; \theta)$
- regression	Linear regression	$\mathcal{N}(h_{\theta}(x), \sigma^2)$
- binary classification	Logistic regression	$\text{Bernoulli}(h_{\theta}(x))$
- multi-class classification	Softmax regression	$\text{Multinomial}([h_{\theta}(x)])$

Can we generalize the linear model to other distributions?



Summary: Linear models

What we've learned so far:

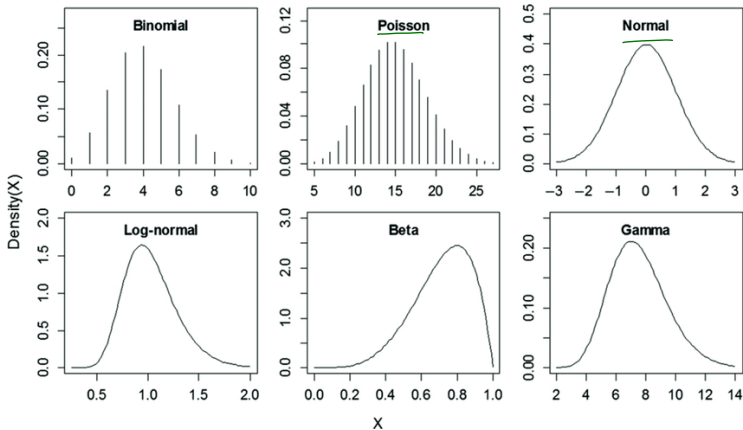
Learning task	Model	$p(y x; \theta)$
regression	Linear regression	$\mathcal{N}(h_\theta(x), \sigma^2)$
binary classification	Logistic regression	Bernoulli($h_\theta(x)$)
multi-class classification	Softmax regression	Multinomial($[h_\theta(x)]$)

Can we generalize the linear model to other distributions?

Generalized Linear Model (GLM): a recipe for constructing linear models in which $p(y|x; \theta)$ is from an **exponential family**.

Review: Exponential Family

Exponential Family of Distributions



Examples of distribution classes in the exponential family.

Exponential Family of Distributions

A class of distributions is in the **exponential family** if its density can be written in the *canonical form*:

$$p(y; \eta) = b(y) e^{\eta^T T(y) - a(\eta)}$$

- ▶ y: random variable
- ▶ η : natural/canonical parameter (that depends on distribution parameter(s)) $\eta = f(\mu)$
- ▶ $T(y)$: sufficient statistic of the distribution
- ▶ $b(y)$: a function of y
- ▶ $a(\eta)$: log partition function (or "cumulant function")

y discrete. $\sum_y p(y; \eta) = 1 \Rightarrow \sum_y \frac{b(y) e^{\eta^T T(y)}}{e^{a(\eta)}} = 1$

$$\frac{1}{e^{a(\eta)}} \sum_y b(y) e^{\eta^T T(y)} = 1$$

$$a(\eta) = \log \left(\sum_y b(y) e^{\eta^T T(y)} \right)$$

Exponential Family

Log partition function $a(\eta)$ is the log of a normalizing constant.
i.e.

$$p(y; \eta) = b(y)e^{\eta^T T(y) - a(\eta)} = \frac{b(y)e^{\eta^T T(y)}}{e^{a(\eta)}}$$

Function $a(\eta)$ is chosen such that $\sum_y p(y; \eta) = 1$
(or $\int_y p(y; \eta) dy = 1$).

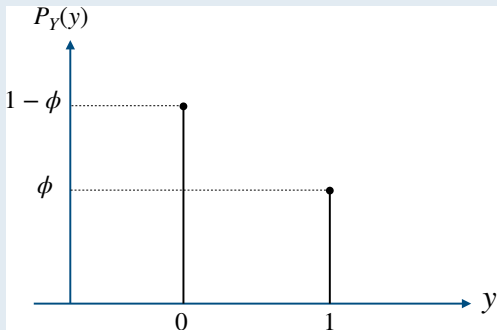
$$a(\eta) = \log \left(\sum_y b(y)e^{\eta^T T(y)} \right)$$

Exponential Family Examples

Bernoulli Distribution

Bernoulli(ϕ): a distribution over $y \in \{0, 1\}$, such that

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$



Bernoulli Distribution

canonical form:

$$p(y; \eta) = \frac{b(y) e^{\eta^T T(y) - a(\eta)}}{Z(\eta)}$$

Bernoulli(ϕ): a distribution over $y \in \{0, 1\}$, such that

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

How to write it in the form of $p(y; \eta) = b(y) e^{\eta^T T(y) - a(\eta)}$?

$$p(y; \phi) = e^{\log p(y; \phi)}$$

$$= e^{y \log \phi + (1-y) \log(1-\phi)}$$

$$= e^{y \log \phi + \log(1-\phi) - y \log(1-\phi)}$$

$$= e^{y \log \frac{\phi}{1-\phi} + \log(1-\phi) - a(\eta)}$$

$$\begin{aligned} & \downarrow \quad \downarrow \quad \downarrow \\ b(y) \quad T(y) \quad \eta \quad \checkmark \\ & \eta = \log \frac{\phi}{1-\phi} \quad \left. \vphantom{\eta} \right\} \text{link} \\ & e^\eta = \frac{\phi}{1-\phi} \end{aligned}$$

$$\begin{aligned} a(\eta) &= -\log(1-\phi) \\ &= -\log\left(1 - \frac{1}{1+e^\eta}\right) \\ &= -\log\left(\frac{1}{1+e^\eta}\right) \\ &= -(-\log(1+e^\eta)) \\ &= \log(1+e^\eta) \quad \checkmark \end{aligned}$$

$$\begin{aligned} e^\eta - e^\eta \phi &= \phi \\ \phi &= \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\eta}} \quad \left. \vphantom{\phi} \right\} \text{response} \end{aligned}$$

(sigmoid function)

Exponential Family Examples

Bernoulli Distribution

Bernoulli(ϕ): a distribution over $y \in \{0, 1\}$, such that

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

- ▶ $\eta =$
- ▶ $b(y) =$
- ▶ $T(y) =$
- ▶ $a(\eta) =$

Exponential Family Examples

Bernoulli Distribution

Bernoulli(ϕ): a distribution over $y \in \{0, 1\}$, such that

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

- ▶ $\eta = \log\left(\frac{\phi}{1-\phi}\right)$
- ▶ $b(y) = 1$
- ▶ $T(y) = y$
- ▶ $a(\eta) = \log(1 + e^\eta)$

Exponential Family Examples

$$p(y; \eta) = \underline{b(y)} e^{\eta^T T(y) - a(\eta)}$$

Gaussian Distribution (unit variance)

Probability density of a Gaussian distribution $\mathcal{N}(\mu, 1)$ over $y \in \mathbb{R}$:

$$\begin{aligned}
 p(y; \theta) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right) \\
 &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y^2 + \mu^2 - 2y\mu)\right) \\
 &= \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2}\right) e^{-\frac{1}{2}(\mu^2 - 2y\mu)} \\
 &= \underbrace{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2}}_{b(y)} e^{\left(y - \frac{\mu^2}{2}\right)} \underbrace{a(\eta)}_{\eta = \mu} = \frac{\mu^2}{2} = \frac{\eta^2}{2}.
 \end{aligned}$$

$\eta = \mu$ $T(y) = y$

Exponential Family Examples

Gaussian Distribution (unit variance)

Probability density of a Gaussian distribution $\mathcal{N}(\mu, 1)$ over $y \in \mathbb{R}$:

$$p(y; \theta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \mu)^2}{2}\right)$$

- ▶ $\eta = \mu$
- ▶ $b(\eta) = \frac{1}{\sqrt{2\pi}} \exp(-\eta^2/2)$
- ▶ $T(y) = y$
- ▶ $a(\eta) = \frac{1}{2}\eta^2$

Exponential Family Examples

Two parameter example:

Gaussian Distribution

Probability density of a Gaussian distribution $\mathcal{N}(\underline{\mu}, \underline{\sigma}^2)$ over $y \in \mathbb{R}$:

$$p(y; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

$$\blacktriangleright \underline{\eta} = \begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix}$$

$$\blacktriangleright b(y) = \frac{1}{\sqrt{2\pi}}$$

$$\blacktriangleright T(y) = \begin{bmatrix} y \\ y^2 \end{bmatrix}$$

$$\blacktriangleright a(\eta) = \frac{\mu^2}{2\sigma^2} + \log \sigma$$

Exponential Family Examples

Poisson distribution: $\text{Poisson}(\lambda)$

Models the probability that an event occurring $y \in \mathbb{N}$ times in a fixed interval of time, *assuming events occur independently at a constant rate*

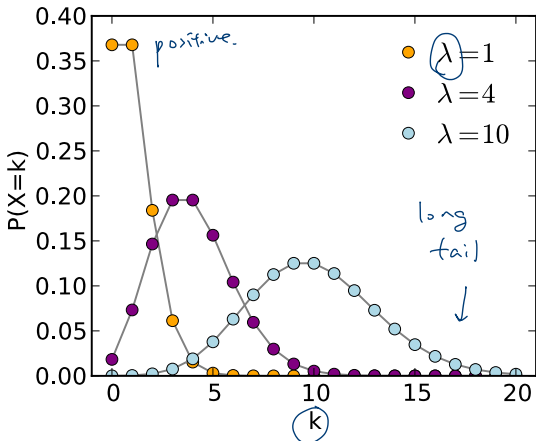
Exponential Family Examples

Poisson distribution: $\text{Poisson}(\lambda)$

Models the probability that an event occurring $y \in \mathbb{N}$ times in a fixed interval of time, *assuming events occur independently at a constant rate*

Probability density function of $\text{Poisson}(\lambda)$ over $y \in \mathcal{Y}$:

$$p(y; \lambda) = \frac{\lambda^y \cdot e^{-\lambda}}{y!}$$



Exponential Family Examples

$$p(y; \eta) = b(y) e^{\eta^T T(y) - a(\eta)}.$$

Poisson distribution $\text{Poisson}(\lambda)$

Probability density function of $\text{Poisson}(\lambda)$ over $y \in \mathcal{Y}$:

$$\begin{aligned}
 p(y; \lambda) &= \frac{1}{y!} e^{\log(\lambda^y e^{-\lambda})} & p(y; \lambda) &= \frac{\lambda^y e^{-\lambda}}{y!} \\
 &= \frac{1}{y!} e^{y \log \lambda - \lambda} & a(\eta) &= \lambda = e^\eta. \\
 &\underbrace{\frac{1}{y!}}_{b(y)} e^{\underbrace{y \log \lambda}_{\eta^T T(y)} - \underbrace{\lambda}_{e^\eta}} & \eta &= \log(\lambda) \\
 & & e^\eta &= \lambda
 \end{aligned}$$

Exponential Family Examples

Poisson distribution $\text{Poisson}(\lambda)$

Probability density function of $\text{Poisson}(\lambda)$ over $y \in \mathcal{Y}$:

$$p(y; \lambda) = \frac{\lambda^y e^{-\lambda}}{y!}$$

- ▶ $\eta = \log \lambda$
- ▶ $b(y) = \frac{1}{y!}$
- ▶ $T(y) = y$
- ▶ $a(\eta) = e^\eta$

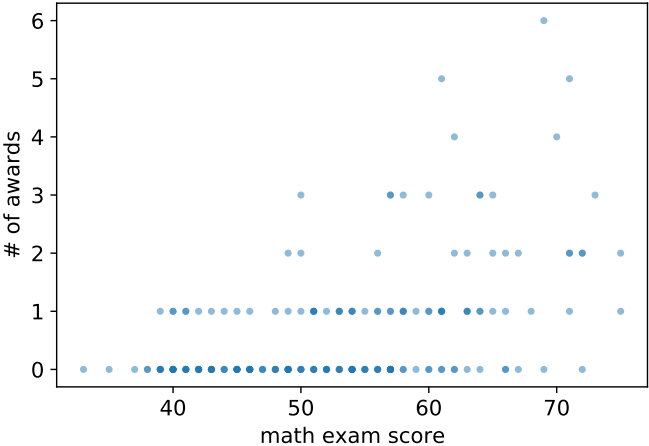
Generalized Linear Models

Generalized Linear Models: Intuition

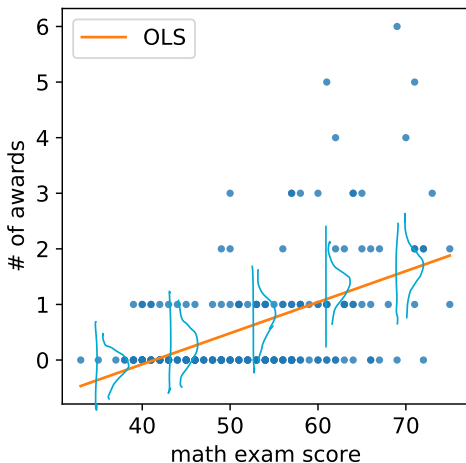
Example 1: Award Prediction

Predict y , **the number of school awards** a student gets given x , the math exam score.

discrete



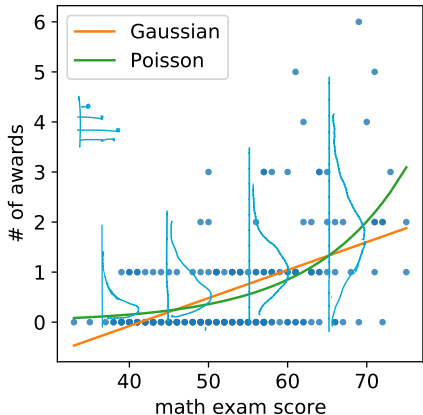
Generalized Linear Models: Intuition



Problems with linear regression:

- ▶ Assumes $y|x; \theta$ has a Normal distribution.
- ▶ Assumes change in x is proportional to change in y

Generalized Linear Models: Intuition



Problems with linear regression:

- ▶ Assumes $y|x; \theta$ has a Normal distribution.
Poisson distribution is better for modeling occurrences
- ▶ Assumes change in x is proportional to change in y
More realistic to be proportional to the rate of increase in y (e.g. doubling or halving y)

Generalized Linear Models : Intuition

Generalized Linear Model (GLM): a recipe for constructing linear models in which $y|x; \theta$ is from an exponential family.

Design motivation of GLM

- ▶ We can select a distribution for **Response variables** y
- ▶ Allow (the **canonical link function of y**) to vary linearly with the input values x

e.g. $\log(\lambda) = \theta^T x$

Nelder, John Ashworth, and Robert William Maclagan Wedderburn. 1972. Generalized Linear Models. Journal of the Royal Statistical Society. Series A (General) 135 (3): 37084.

Generalized Linear Models: Construction

Formal GLM assumptions & design decisions:

1. $y|x; \theta \sim \text{ExponentialFamily}(\eta)$
e.g. Gaussian, Poisson, Bernoulli, Multinomial, Beta ...
2. The hypothesis function $h(x)$ is $\mathbb{E}[T(y)|x]$,
e.g. When $T(y) = y$, $h(x) = \mathbb{E}[y|x]$ *sufficient statistics of y*
3. The natural parameter η and the inputs x are related linearly:

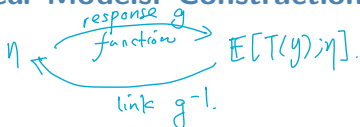
η is a number:

$$\eta = \theta^T x$$

η is a vector:

$$\eta_i = \theta_i^T x \quad \forall i = 1, \dots, n \quad \text{or} \quad \eta = \Theta^T x$$

Generalized Linear Models: Construction



Relate natural parameter η to distribution mean $\mathbb{E}[T(y); \eta]$:

- ▶ **Canonical response function** g gives the mean of the distribution

$$g(\eta) = \mathbb{E}[T(y); \eta]$$

a.k.a. the “mean function”

Generalized Linear Models: Construction

Relate natural parameter η to distribution mean $\mathbb{E}[T(y); \eta]$:

- ▶ **Canonical response function** g gives the mean of the distribution

$$g(\eta) = \mathbb{E}[T(y); \eta]$$

a.k.a. the “mean function”

- ▶ g^{-1} is called the **canonical link function**

$$\eta = g^{-1}(\mathbb{E}[T(y); \eta])$$

GLM example: ordinary least square

Apply GLM construction rules:

1. Let $y|x; \theta \sim N(\mu, 1)$

$$\underline{\eta = \mu}, \quad \underline{T(y) = y}$$

GLM example: ordinary least square

Apply GLM construction rules:

1. Let $y|x; \theta \sim N(\mu, 1)$

$$\underline{\eta} = \mu, \quad \underline{T(y)} = \underline{y}$$

2. Derive hypothesis function:

$$\begin{aligned} \underline{h_{\theta}(x)} &= \underline{\mathbb{E}[T(y)|x; \theta]} \\ &= \underline{\mathbb{E}[y|x; \theta]} \quad y|x; \theta \sim N(\mu, 1) \\ &= \underline{\mu} = \underline{\eta} \end{aligned}$$

GLM example: ordinary least square

Apply GLM construction rules:

1. Let $y|x; \theta \sim N(\mu, 1)$

$$\eta = \mu, T(y) = y$$

2. Derive hypothesis function:

$$\begin{aligned}h_{\theta}(x) &= \mathbb{E}[T(y)|x; \theta] \\ &= \mathbb{E}[y|x; \theta] \\ &= \mu = \eta\end{aligned}$$

3. Adopt linear model $\eta = \theta^T x$:

$$\underline{h_{\theta}(x) = \eta = \theta^T x}$$

GLM example: ordinary least square

Apply GLM construction rules:

1. Let $y|x; \theta \sim N(\mu, 1)$

$$\eta = \mu, T(y) = y$$

2. Derive hypothesis function:

$$\begin{aligned}h_{\theta}(x) &= \mathbb{E}[T(y)|x; \theta] \\ &= \mathbb{E}[y|x; \theta] \\ &= \mu = \eta\end{aligned}$$

3. Adopt linear model $\eta = \theta^T x$:

$$h_{\theta}(x) = \eta = \theta^T x$$

Canonical response function: $\underline{\mu} = g(\underline{\eta}) = \underline{\eta}$ (identity)

Canonical link function: $\underline{\eta} = g^{-1}(\underline{\mu}) = \underline{\mu}$ (identity)

GLM example: logistic regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Bernoulli}(\phi)$ $g^{-1}(\phi)$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right), \quad T(y) = y$$

GLM example: logistic regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Bernoulli}(\phi)$ $g^{-1}(\phi)$
 $\eta = \log\left(\frac{\phi}{1-\phi}\right)$, $T(y) = y$
2. Derive hypothesis function:

$$\begin{aligned}h_{\theta}(x) &= \mathbb{E}[T(y)|x; \theta] \\&= \mathbb{E}[y|x; \theta] \quad \downarrow g(\eta) \\&= \underline{\underline{\phi}} = \underline{\underline{\frac{1}{1 + e^{-\eta}}}}\end{aligned}$$

GLM example: logistic regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Bernoulli}(\phi)$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right), \quad T(y) = y$$

2. Derive hypothesis function:

$$\begin{aligned} h_{\theta}(x) &= \mathbb{E}[T(y)|x; \theta] \\ &= \mathbb{E}[y|x; \theta] \\ &= \phi = \frac{1}{1 + e^{-\eta}} \end{aligned}$$

3. Adopt linear model $\eta = \underline{\theta^T x}$:

$$h_{\theta}(x) = \underline{\frac{1}{1 + e^{-\theta^T x}}}$$

GLM example: logistic regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Bernoulli}(\phi)$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right), \quad T(y) = y$$

2. Derive hypothesis function:

$$\begin{aligned} h_{\theta}(x) &= \mathbb{E}[T(y)|x; \theta] \\ &= \mathbb{E}[y|x; \theta] \\ &= \phi = \frac{1}{1 + e^{-\eta}} \end{aligned}$$

3. Adopt linear model $\eta = \theta^T x$:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Canonical response function: $\phi = g(\eta) = \text{sigmoid}(\eta)$

GLM example: logistic regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Bernoulli}(\phi)$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right), \quad T(y) = y$$

logit(ϕ)
↓

2. Derive hypothesis function:

$$\begin{aligned} h_{\theta}(x) &= \mathbb{E}[T(y)|x; \theta] \\ &= \mathbb{E}[y|x; \theta] \\ &= \phi = \frac{1}{1 + e^{-\eta}} \end{aligned}$$

sigmoid(η)

3. Adopt linear model $\eta = \theta^T x$:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Canonical response function: $\phi = g(\eta) = \text{sigmoid}(\eta)$

Canonical link function : $\eta = g^{-1}(\phi) = \text{logit}(\phi)$

GLM example: Poisson regression

Example 1: Award Prediction

Predict y , **the number of school awards** a student gets given x , the math exam score.

Use GLM to find the hypothesis function...

GLM example: Poisson regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Poisson}(\lambda)$

$$\eta = \log(\lambda), \quad T(y) = y$$

2. Derive hypothesis function:

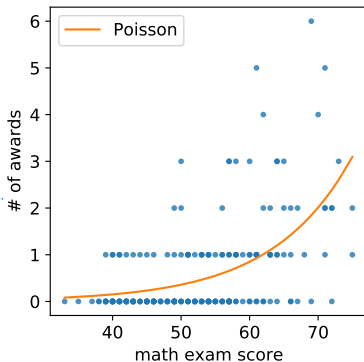
$$h_{\theta}(x) = \mathbb{E}[T(y)|x; \theta]$$
$$= \lambda = e^{\eta} \quad \text{response fun.}$$

3. Adopt linear model $\eta = \theta^T x$:

$$h_{\theta}(x) = e^{\theta^T x}$$

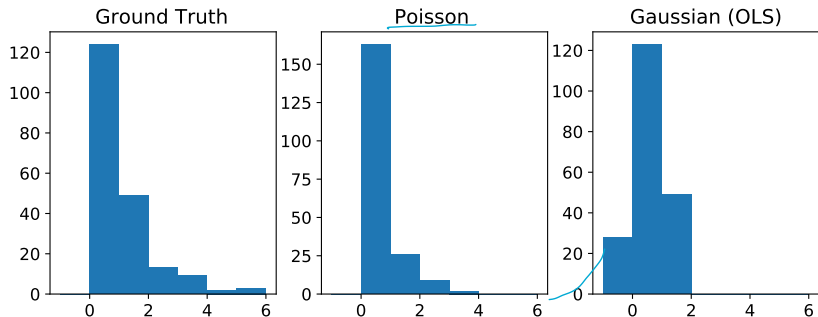
Canonical response function: $\lambda = g(\eta) = e^{\eta}$

Canonical link function: $\eta = g^{-1}(\lambda) = \log(\lambda)$



GLM example: Poisson regression

Distribution of the predicted number of awards (y)



Poisson regression successfully captures the long tail of $P(y)$

GLM example: Softmax regression

$$P(y|x) \sim \text{Multinomial}(\underbrace{(\phi_1, \dots, \phi_k)}_{k-1})$$

Probability mass function of a Multinomial distribution over k outcomes

$$P(y; \eta) = b(y) e^{\eta^T T(y) - a(\eta)}$$

$$P(y; \phi) = \prod_{i=1}^k \phi_i^{1\{y=i\}} \begin{cases} 1 & y=i \\ 0 & y \neq i \end{cases}$$

Derive the exponential family form of Multinomial(ϕ_1, \dots, ϕ_k): **Note:**

$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$ is not a parameter

Let $\partial(y)_i = 1\{y=i\} \in \{0,1\}$.

$$P(y; \phi) = \left(\prod_{i=1}^{k-1} \phi_i^{1\{y=i\}} \right) \cdot \phi_k^{1\{y=k\}}$$

$$\partial(y) = \begin{bmatrix} \partial(y)_1 \\ \vdots \\ \partial(y)_k \end{bmatrix} = \begin{bmatrix} 1\{y=1\} \\ \vdots \\ 1\{y=k\} \end{bmatrix}$$

$\partial(y)_k = 1 - \sum_{i=1}^{k-1} \partial(y)_i$

$$\sum_{i=1}^{k-1} \partial(y)_i = 1$$

$$= e^{\sum_{i=1}^{k-1} \partial(y)_i \log \phi_i + \partial(y)_k \log \phi_k}$$

$$= e^{\sum_{i=1}^{k-1} \partial(y)_i \log \phi_i + \log \phi_k - \sum_{i=1}^{k-1} \partial(y)_i \log \phi_k}$$

$$a(\eta) = -\log \phi_k = -\log \left(\frac{1}{\sum_{i=1}^k e^{\eta_i}} \right) = \log \sum_{i=1}^k e^{\eta_i}$$

$$= e^{\sum_{i=1}^{k-1} \partial(y)_i \log \frac{\phi_i}{\phi_k} + \log \phi_k}$$

$$\rightarrow \begin{bmatrix} \partial(y)_1 \\ \vdots \\ \partial(y)_k \end{bmatrix}^T \begin{bmatrix} \log \frac{\phi_1}{\phi_k} \\ \vdots \\ \log \frac{\phi_{k-1}}{\phi_k} \end{bmatrix}$$

$T(y) = \partial(y)$

$$\eta_i = \log \frac{\phi_i}{\phi_k}$$

$$e^{\eta_i} = \frac{\phi_i}{\phi_k}$$

$$\phi_i = \phi_k \cdot e^{\eta_i}$$

$$= \frac{e^{\eta_i}}{\sum_{i=1}^k e^{\eta_i}} \} g(\eta)$$

since $\sum_{i=1}^k \phi_k e^{\eta_i} = 1$

$$\phi_k = \frac{1}{\sum_{i=1}^k e^{\eta_i}}$$

GLM example: Softmax regression

Probability mass function of a Multinomial distribution over k outcomes

$$p(y; \phi) = \prod_{i=1}^k \phi_i^{\mathbf{1}\{y=i\}}$$

Derive the exponential family form of Multinomial(ϕ_1, \dots, ϕ_k): **Note:**

$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$ is not a parameter

$\eta(y)_1$

$\eta(y)_k$

$$\underline{T(y)} = \begin{bmatrix} \mathbf{1}\{y=1\} \\ \vdots \\ \mathbf{1}\{y=k-1\} \end{bmatrix}$$

$$T(y)_i = \mathbf{1}\{y=i\} = \begin{cases} 0 & y \neq i \\ 1 & y = i \end{cases}$$

GLM example: Softmax regression

Probability mass function of a Multinomial distribution over k outcomes

$$p(y; \phi) = \prod_{i=1}^k \phi_i^{\mathbf{1}\{y=i\}}$$

Derive the exponential family form of Multinomial(ϕ_1, \dots, ϕ_k): **Note:**

$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$ is not a parameter

GLM example: Softmax regression

Probability mass function of a Multinomial distribution over k outcomes

$$p(y; \phi) = \prod_{i=1}^k \phi_i^{\mathbf{1}\{y=i\}}$$

Derive the exponential family form of Multinomial(ϕ_1, \dots, ϕ_k): **Note:**

$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$ is not a parameter

$$\blacktriangleright T(y) = \begin{bmatrix} \mathbf{1}\{y=1\} \\ \vdots \\ \mathbf{1}\{y=k-1\} \end{bmatrix}$$

$$T(y)_i = \mathbf{1}\{y=i\} = \begin{cases} 0 & y \neq i \\ 1 & y = i \end{cases}$$

$$\blacktriangleright a(\eta) = -\log(\phi_k)$$

$$\blacktriangleright \eta = \begin{bmatrix} \log\left(\frac{\phi_1}{\phi_k}\right) \\ \vdots \\ \log\left(\frac{\phi_{k-1}}{\phi_k}\right) \end{bmatrix}$$

$$\blacktriangleright \underline{b(y)} = 1$$

GLM example: Softmax regression

Apply GLM construction rules:

1. Let $\underline{y|x}; \theta \sim \text{Multinomial}(\phi_1, \dots, \phi_k)$, for all $i = 1 \dots k - 1$

$$\underline{\eta_i} = \log\left(\frac{\phi_i}{\phi_k}\right), \quad \underline{T(y)} = \begin{bmatrix} \mathbf{1}\{y = 1\} \\ \vdots \\ \mathbf{1}\{y = k - 1\} \end{bmatrix}$$

GLM example: Softmax regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Multinomial}(\phi_1, \dots, \phi_k)$, for all $i = 1 \dots k - 1$

$$\eta_i = \log\left(\frac{\phi_i}{\phi_k}\right), \quad T(y) = \begin{bmatrix} \mathbf{1}\{y = 1\} \\ \vdots \\ \mathbf{1}\{y = k - 1\} \end{bmatrix}$$

Compute inverse: $\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} \leftarrow \text{canonical response function}$

GLM example: Softmax regression

Apply GLM construction rules:

1. Let $y|x; \theta \sim \text{Multinomial}(\phi_1, \dots, \phi_k)$, for all $i = 1 \dots k - 1$

$$\eta_i = \log\left(\frac{\phi_i}{\phi_k}\right), \quad \underline{T(y)} = \begin{bmatrix} \mathbf{1}\{y = 1\} \\ \vdots \\ \mathbf{1}\{y = k - 1\} \end{bmatrix}$$

Compute inverse: $\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} \leftarrow \text{canonical response function}$

2. Derive hypothesis function:

$$h_{\theta}(x) = \mathbb{E} \left[\begin{bmatrix} \mathbf{1}\{y = 1\} \\ \vdots \\ \mathbf{1}\{y = k - 1\} \end{bmatrix} \middle| x; \theta \right] = \begin{bmatrix} \phi_1 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$$\underline{\phi_i} = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} \quad \theta_i^T x$$

GLM example: Softmax regression

3. Adopt linear model $\eta_i = \theta_i^T x$:

$$\phi_i = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}} \text{ for all } i = 1 \dots k - 1$$

$$h_{\theta}(x) = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_{k-1}^T x} \end{bmatrix}$$

GLM example: Softmax regression

3. Adopt linear model $\eta_i = \theta_i^T x$:

$$\phi_i = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}} \text{ for all } i = 1 \dots k - 1$$

$$h_{\theta}(x) = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x}} \begin{bmatrix} e^{\theta_1^T x} \\ \vdots \\ e^{\theta_{k-1}^T x} \end{bmatrix}$$

Canonical response function: $\phi_i = g(\eta) = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$ ✓

Canonical link function : $\eta_i = g^{-1}(\phi_i) = \log\left(\frac{\phi_i}{\phi_k}\right)$ ✓

GLM Summary

Sufficient statistic $T(y)$

Response function $g(\eta)$

Link function $g^{-1}(\mathbb{E}[T(y); \eta])$

$$\mu \sim \theta^T X$$

$$\log \frac{\phi}{1-\phi} \sim \theta^T X$$

$$\log(\lambda) \sim \theta^T X$$

Exponential Family	\mathcal{Y}	$T(y)$	$g(\eta)$	$g^{-1}(\mathbb{E}[T(y); \eta])$
$\mathcal{N}(\mu, 1)$	\mathbb{R}	y	η	μ
Bernoulli(ϕ)	$\{0, 1\}$	y	$\frac{1}{1+e^{-\eta}}$	$\log \frac{\phi}{1-\phi}$
Poisson(λ)	\mathbb{N}	y	e^η	$\log(\lambda)$
Multinomial(ϕ_1, \dots, ϕ_k)	$\{1, \dots, k\}$	$\mathbf{1}\{y = i\}$	$\frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$	$\eta_i = \log\left(\frac{\phi_i}{\phi_k}\right)$

GLM is effective for modelling different types of distributions over y