# Generalized Linear Model

Mohammad Emtiyaz Khan
EPFL
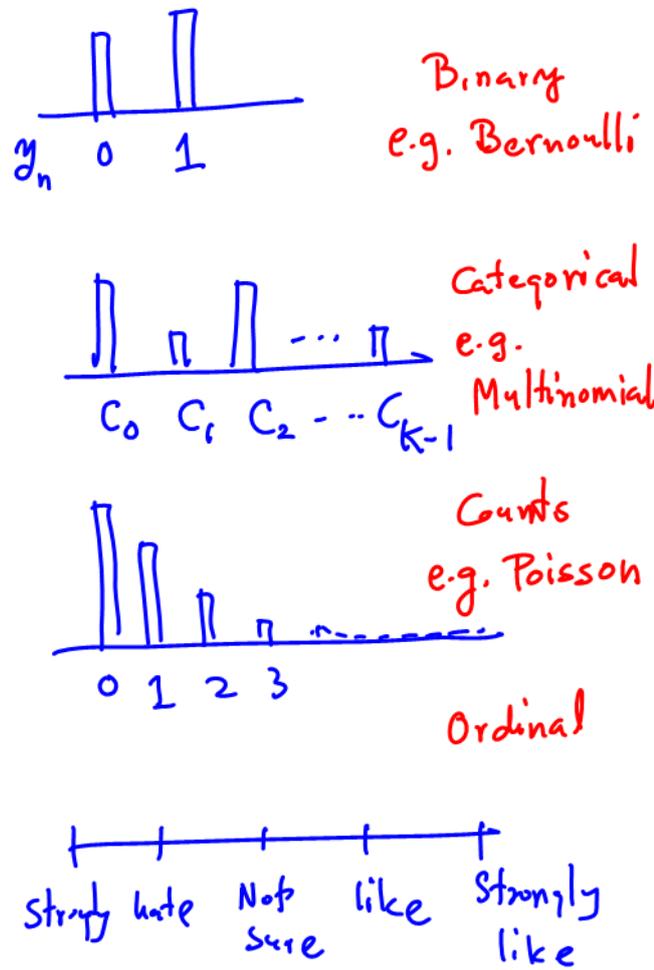
Oct 15, 2015

# Motivation

Logistic regression extends linear regression to binary outputs. Can we generalize this procedure to other kind of outputs, e.g. categorical, ordinal, counts, etc?

Specifically, can we generalize the maximum-likelihood procedure to a generic distribution $p(y_n|\widetilde{\mathbf{x}}_n^T\boldsymbol{\beta})$?

The answer is yes, we can do this for exponential family distributions.

*(margin annotations, red)* Binary e.g. Bernoulli

$y_n$  0  1

Categorical e.g. Multinomial

$C_0$  $C_1$  $C_2$ --- $C_{K-1}$

Counts e.g. Poisson

0 1 2 3

Ordinal

*(margin, blue)* Strongly hate   Not Sure   like   Strongly like

# Logistic regression revisited

In logistic regression, we used the following distribution:

$\eta_n = \widetilde{\mathbf{x}}_n^T\boldsymbol{\beta}$

$$p(y_n|\eta_n) := \frac{\exp(y_n\eta_n)}{1 + e^{\eta_n}} = \frac{e^{y_n\eta_n}}{e^{\log(1+e^{\eta_n})}} = e^{y_n\eta_n - \log(1+e^{\eta_n})}$$

$$= \exp\left[\underbrace{y_n\eta_n}_{linear} - \underbrace{\log(1 + e^{\eta_n})}_{A(\eta_n)}\right].$$

This specific form of the distribution is not coincidental.

In general, we require the following form of the distribution:

$$p(y_n|\eta_n) := \exp\left[\underbrace{y_n\eta_n}_{linear} - \underbrace{A(\eta_n)}_{function\ of\ \eta_n}\right] h(y_n),$$

where $h$ does not depend on $\eta_n$.

1

# Exponential family distribution

Exponential family distribution is a general class of distributions defined as shown below,

$$Z = \int h(y) \exp\left[\underline{\eta}^T \underline{\phi}(y) - A(\eta)\right] dy$$

natural parameter

log-partition function

$$p(y|\boldsymbol{\eta}) := \frac{\overbrace{h(y)}}{Z} \exp\left[\boldsymbol{\eta}^T \boldsymbol{\phi}(y) - A(\boldsymbol{\eta})\right].$$

Partition function

$\longrightarrow$ Sufficient statistics

This definition is the *canonical* form and assumes a scalar-valued random variable $y$. See Wikipedia for a more general definition.

Bernoulli distribution

Example: <u>Bernoulli</u> distribution is in exponential family.

$$p\left(y=1/\mu\right) = \mu$$
$$p\left(y=0/\mu\right) = 1-\mu,$$

$$\mu \in (0,1)$$

$$p(y|\mu) := \mu^y(1-\mu)^{1-y}, \text{ where } \mu \in (0,1)$$

$$= \exp\left[y\log\frac{\mu}{1-\mu} + \log(1-\mu)\right]$$

$\phi(y)$    $\eta$    $-A(\eta)$

$$\mu^y(1-\mu)^{1-y}$$

$$= \exp\left[\log\left\{\mu^y(1-\mu)^{1-y}\right\}\right]$$

There is a relationship between $\eta$ and $\mu$ through the link function,

$$= \exp\left[y\log\mu + (1-y)\log(1-\mu)\right]$$

$$\eta = \log\frac{\mu}{1-\mu} \iff \mu = \frac{e^\eta}{1+e^\eta}.$$

$g(\mu)$

Note that $\boxed{\mu = \mathbb{E}(\phi(y))}$ is the mean parameter of $y$ (and $y$ is the *sufficient* statistics to describe the distribution).
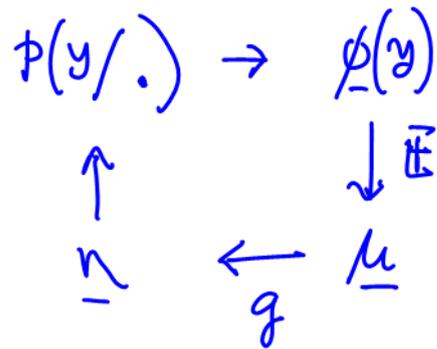
$$A(\eta) = -\log(1-\mu)$$
$$= -\log\left(1 - \frac{e^\eta}{1+e^\eta}\right)$$
$$= +\log(1+e^\eta)$$

# Useful properties of exponential family

In general, there is a relationship between the mean $\boldsymbol{\mu} := \mathbb{E}(\boldsymbol{\phi}(y))$ and $\boldsymbol{\eta}$ defined using a link function $g$.

$$\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\mu}) \iff \boldsymbol{\mu} = \mathbf{g}^{-1}(\boldsymbol{\eta}).$$

See the table of link functions and many other examples of exponential family in the KPM book chapter on "Generalized Linear Model".

$$P(y/\cdot) \rightarrow \underline{\phi}(y)$$

$$\uparrow \qquad \downarrow \mathbb{E}$$

$$\underline{\eta} \xleftarrow{g} \underline{\mu}$$

Derive this for Gaussian and multinomial distribution

Another useful property is that the first and second derivatives of $A(\eta)$ are related to the mean and the variance of the sufficient statistics.

$$\frac{\partial A(\eta)}{\partial \eta} = \mathbb{E}[\underline{\phi}(\eta)] \ , \ \frac{\partial^2 A(\eta)}{\partial \eta^2} = \mathrm{Var}[\underline{\phi}(\eta)].$$

Therefore, $A$ is convex.

Using this property, we can generalize the normal equation.

# The maximum likelihood estimate

The generalized maximum likelihood cost to minimize is,

$$\min \mathcal{L}(\boldsymbol{\beta}) := -\sum_{n=1}^{N} \log p(y_n | \widetilde{\mathbf{x}}_n^T \boldsymbol{\beta}),$$

where $p(y_n | \eta_n)$ is an exponential family distribution with $\eta_n = \widetilde{\mathbf{x}}_n^T \boldsymbol{\beta}$ as the natural parameters.

The maximum likelihood solution can be obtained by solving the following normal equation:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\beta}} = \widetilde{\mathbf{X}}^T \left[ \mathbf{g}^{-1}(\boldsymbol{\eta}) - \underline{\phi}(\underline{y}) \right]$$

where $\boldsymbol{\eta}$ is a vector of all $\eta_n$ and the function $\mathbf{g}^{-1}(\boldsymbol{\eta})$ is a vector of $\mathbf{g}^{-1}(\eta_n)$.

**To do**

1. Read the following sections in the KPM book: Section 9.1.1 to 9.1.4, Section 9.3.1 to 9.3.2.

2. Derive exponential family form for Gaussian distribution and multinomial distributions.

3. Derive generalized linear model for regression with Poisson distribution for count data.

*Handwritten annotations:*

$p\left(y_n / \eta_n\right)$ — Find a $\underline{\beta}$ that is a "good-fit" to my data $\left(y_n, \underline{x}_n\right)$

$\eta_n := \underline{x}_n^T \underline{\beta}$

Assume $\eta$ is a scalar

$\mathcal{L}_n(\underline{\beta}) := -\log p\left(y_n / \underline{x}_n^T \underline{\beta}\right)$

$= -\sum_{n=1}^{N} \left[ \eta_n \phi(y) - A(\eta_n) \right] + cnst$

$\frac{\partial \mathcal{L}_n}{\partial \eta_n} = -\phi(y_n) + \frac{\partial}{\partial \eta_n} A(\eta_n)$

$= -\phi(y_n) + \mu_n$

$= -\phi(y_n) + \bar{g}^{-1}(\eta_n)$

$\frac{\partial \mathcal{L}}{\partial \underline{\beta}} = \widetilde{X}^T \left[ \bar{g}^{-1}(\eta) - \phi(\underline{y}) \right]$ ← Vectorize

$\frac{\partial}{\partial \underline{\beta}} = \frac{\partial \eta_n}{\partial \underline{\beta}} \frac{\partial \mathcal{L}_n}{\partial \eta_n} = \widetilde{x}_n \frac{\partial \mathcal{L}_n}{\partial \eta_n}$

Hessian

$\frac{\partial^2 \mathcal{L}}{\partial \underline{\beta} \partial \underline{\beta}^T} = X^T S X$, where $S$ is a diagonal matrix with $S_{nn} = \frac{\partial^2 A(\eta_n)}{\partial \eta_n}$ with $\eta_n = \widetilde{X}_n^T \underline{\beta}$

↳ always pos. semi-definite