# The Exponential Family & Generalized Linear Models

## CSC2541 Tutorial 2, Winter 2022

Jenny Bao

Jan 20, 2022

# Overview

The Exponential Family
- ▶ Formula & basics
- ▶ Examples: Bernoulli, Gaussian, ...
- ▶ Useful identities
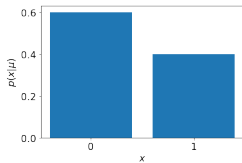
Generalized Linear Models

# The Exponential Family

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

- ▶ $\boldsymbol{\eta}$: natural parameters
- ▶ $\boldsymbol{u}(\boldsymbol{x})$: sufficient statistic
- ▶ $\mathcal{Z}(\boldsymbol{\eta})$: partition function, ensures the distribution $p(\boldsymbol{x}|\boldsymbol{\eta})$ is normalized.

$$\text{continuous:} \quad \mathcal{Z}(\boldsymbol{\eta}) = \int h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x})) d\boldsymbol{x}$$

$$\text{discrete:} \quad \mathcal{Z}(\boldsymbol{\eta}) = \sum_{\boldsymbol{x}} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

# Example 1: Bernoulli



$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

Bernoulli distribution:
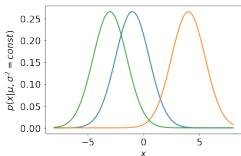
$$p(x|\mu) = \mathrm{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x}$$

Put in the exponential family form:

$$\mathrm{Bern}(x|\mu) = (1-\mu)\big(\frac{\mu}{1-\mu}\big)^x$$

$$= \underbrace{(1-\mu)}_{\frac{1}{\mathcal{Z}(\eta)}} \cdot \underbrace{1}_{h(x)} \cdot \exp\big\{ \underbrace{\big(\log \frac{\mu}{1-\mu}\big)}_{\eta} \cdot \underbrace{x}_{u(x)} \big\}$$

▶ $\implies \mu = \sigma(\eta),\ \mathcal{Z}(\eta) = \sigma(-\eta)$

Exercise: put the multinomial distribution in the standard form for exponential family
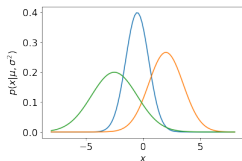
# Example 2: Gaussian ($\mu$)



$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

Gaussian distribution (treating only $\mu$ as parameter, assuming $\sigma$ is constant):

$$
\begin{aligned}
p(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \\
&= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2}\right) \\
&= \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right)}_{h(x)} \cdot \underbrace{\exp\left(-\frac{\mu^2}{2\sigma^2}\right)}_{\frac{1}{\mathcal{Z}(\boldsymbol{\eta})}} \exp\left(\underbrace{\frac{\mu}{\sigma^2}}_{\boldsymbol{\eta}^\top} \underbrace{x}_{u(x)}\right)
\end{aligned}
$$

# Example 3: Gaussian ($\mu$ and $\sigma$)



$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

Gaussian distribution (treating both $\mu$ and $\sigma$ as parameters):

$$
\begin{aligned}
p(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}} \sigma^{-1} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \\
&= \frac{1}{\sqrt{2\pi}} \sigma^{-1} \exp\left(-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2}\right) \\
&= \underbrace{\frac{1}{\sqrt{2\pi}}}_{h(x)} \cdot \underbrace{\sigma^{-1}\exp\left(-\frac{\mu^2}{2\sigma^2}\right)}_{\frac{1}{\mathcal{Z}(\eta)}} \exp\left(\underbrace{\begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix}}_{\eta^\top}{}^\top \underbrace{\begin{bmatrix} x \\ x^2 \end{bmatrix}}_{u(x)}\right)
\end{aligned}
$$

# The Exponential Family

Other members of the exponential family:

- Poisson, gamma, exponential, beta, Dirichlet, ...

# Why studying the exponential family?

Many convenient properties
- Sufficient statistics for maximum likelihood
- Many convenient identities for $\mathcal{Z}(\boldsymbol{\eta})$ (the partition function)
  - Relates concepts such as the Fisher information matrix

Can be used to derive the Generalized Linear Models (GLM)

# Maximum likelihood & sufficient statistics

Consider i.i.d. data $\boldsymbol{X} = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$. Find $\boldsymbol{\eta}$ to maximize $p(\boldsymbol{X}|\boldsymbol{\eta})$ (maximum likelihood).

$$p(\boldsymbol{X}|\boldsymbol{\eta}) = \prod_{i=1}^{N} \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} \left( h(\boldsymbol{x}_i) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}_i)) \right)$$

$$= \left( \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} \right)^N \left( \prod_{i=1}^{N} h(\boldsymbol{x}_i) \right) \exp(\boldsymbol{\eta}^\top \sum_{i=1}^{N} \boldsymbol{u}(\boldsymbol{x}_i))$$

Take derivative of the log-likelihood and set it to 0.

$$\nabla_{\boldsymbol{\eta}} \log p(\boldsymbol{X}|\boldsymbol{\eta}) = -N \nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) + \sum_{i=1}^{N} \boldsymbol{u}(\boldsymbol{x}_i) = 0$$

$$\implies \nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{u}(\boldsymbol{x}_i)$$

# Maximum likelihood & sufficient statistics

$$\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{u}(\boldsymbol{x}_i)$$

The maximum-likelihood solution $\boldsymbol{\eta}$ only depends on $\sum_{i=1}^{N} \boldsymbol{u}(\boldsymbol{x}_i)$.

▶ Hence $\boldsymbol{u}(\boldsymbol{x})$ is called the sufficient statistic

Examples:

▶ **Bernoulli:** $u(x) = x$. Only need to store $\sum_i x_i$.

▶ **Gaussian ($\mu$ and $\sigma$):** $u(x) = \begin{bmatrix} x & x^2 \end{bmatrix}^{\top}$. Need to store both $\sum_i x_i$ and $\sum_i x_i^2$.

# Identities

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

1. $\boxed{\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})] =: \boldsymbol{\xi}}$ (moments)

Derivation:

$$\begin{aligned}
\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) &= \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} \nabla_{\boldsymbol{\eta}} \int h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x})) d\boldsymbol{x} \\
&= \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} \int h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x})) \boldsymbol{u}(\boldsymbol{x}) d\boldsymbol{x} \\
&= \int \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x})) \boldsymbol{u}(\boldsymbol{x}) d\boldsymbol{x} \\
&= \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})]
\end{aligned}$$

# Identities

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

1. $\boxed{\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})] =: \boldsymbol{\xi}}$ (moments)

- ▶ There's a 1-to-1 mapping between $\boldsymbol{\eta} \leftrightarrow \boldsymbol{\xi}$
- ▶ $\boldsymbol{\xi}$ is an alternative parameterization for the exponential family

$$p(\boldsymbol{x}|\boldsymbol{\eta}) \leftrightarrow p(\boldsymbol{x}|\boldsymbol{\xi})$$

# Identities

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

1. $\boxed{\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})] =: \boldsymbol{\xi}}$ (moments)

2. $\boxed{\nabla_{\boldsymbol{\eta}} \log p(\boldsymbol{x}|\boldsymbol{\eta}) = \boldsymbol{u}(\boldsymbol{x}) - \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})]}$

Derivation:

$$\begin{aligned}
\nabla_{\boldsymbol{\eta}} \log p(\boldsymbol{x}|\boldsymbol{\eta}) &= -\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) + \nabla_{\boldsymbol{\eta}}(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x})) \\
&= -\mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})] + \boldsymbol{u}(\boldsymbol{x})
\end{aligned}$$

# Identities

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

1. $\boxed{\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{\eta})}[\boldsymbol{u}(\boldsymbol{x})] =: \boldsymbol{\xi}}$ (moments)

2. $\boxed{\nabla_{\boldsymbol{\eta}} \log p(\boldsymbol{x}|\boldsymbol{\eta}) = \boldsymbol{u}(\boldsymbol{x}) - \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{\eta})}[\boldsymbol{u}(\boldsymbol{x})]}$

Recall, in maximum likelihood, we have

$$\nabla_{\boldsymbol{\eta}} \sum_{i=1}^{N} \log p(\boldsymbol{x}_i|\boldsymbol{\eta}) = N\left( \underbrace{\frac{1}{N} \sum_{i=1}^{N} \boldsymbol{u}(\boldsymbol{x}_i)}_{\text{empirical moments } \hat{\boldsymbol{\xi}}} - \underbrace{\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta})}_{\text{moments } \boldsymbol{\xi}} \right)$$

$$= N \cdot (\hat{\boldsymbol{\xi}} - \boldsymbol{\xi}) = 0$$

Maximum likelihood $\rightarrow$ moment matching

# Identities

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

1. $\boxed{\nabla_{\boldsymbol{\eta}} \log \mathcal{Z}(\boldsymbol{\eta}) = \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})] =: \boldsymbol{\xi}}$ (moments)

2. $\boxed{\nabla_{\boldsymbol{\eta}} \log p(\boldsymbol{x}|\boldsymbol{\eta}) = \boldsymbol{u}(\boldsymbol{x}) - \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})]}$

3.1. $\boxed{\nabla_{\boldsymbol{\eta}}^2 \log \mathcal{Z}(\boldsymbol{\eta}) = \mathrm{Cov}(\boldsymbol{u}(\boldsymbol{x})) = -\nabla_{\boldsymbol{\eta}}^2 \log p(\boldsymbol{x}|\boldsymbol{\eta})}$

Derivation: $\nabla_{\boldsymbol{\eta}}^2 \log p(\boldsymbol{x}|\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \boldsymbol{u}(\boldsymbol{x}) - \nabla_{\boldsymbol{\eta}}^2 \log \mathcal{Z}(\boldsymbol{\eta}) = -\nabla_{\boldsymbol{\eta}}^2 \log \mathcal{Z}(\boldsymbol{\eta})$

$$\nabla_{\boldsymbol{\eta}}^2 \log \mathcal{Z}(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})]$$

the grad-log trick $\rightarrow$ $\quad = \mathbb{E}_{x \sim p(x|\eta)}[\nabla_{\boldsymbol{\eta}} \log p(\boldsymbol{x}|\boldsymbol{\eta}) \boldsymbol{u}(\boldsymbol{x})^\top]$

$\quad = \mathbb{E}_{x \sim p(x|\eta)}[(\boldsymbol{u}(\boldsymbol{x}) - \mathbb{E}_{x \sim p(x|\eta)}[\boldsymbol{u}(\boldsymbol{x})]) \boldsymbol{u}(\boldsymbol{x})^\top]$

$\mathbb{E}[(\boldsymbol{u} - \mathbb{E}[\boldsymbol{u}])\mathbb{E}[\boldsymbol{u}]] = 0 \rightarrow$ $\quad = \mathbb{E}[(\boldsymbol{u}(\boldsymbol{x}) - \mathbb{E}[\boldsymbol{u}(\boldsymbol{x})])(\boldsymbol{u}(\boldsymbol{x}) - \mathbb{E}[\boldsymbol{u}(\boldsymbol{x})])^\top]$

$\quad = \mathrm{Cov}(\boldsymbol{u}(\boldsymbol{x}))$

# Identities

$$p(x|\eta) = \frac{1}{\mathcal{Z}(\eta)} h(x) \exp(\eta^\top u(x))$$

1. $\boxed{\nabla_\eta \log \mathcal{Z}(\eta) = \mathbb{E}_{x \sim p(x|\eta)}[u(x)] =: \xi}$ (moments)

2. $\boxed{\nabla_\eta \log p(x|\eta) = u(x) - \mathbb{E}_{x \sim p(x|\eta)}[u(x)]}$

3.1. $\boxed{\nabla_\eta^2 \log \mathcal{Z}(\eta) = \mathrm{Cov}(u(x)) = -\nabla_\eta^2 \log p(x|\eta)}$

3.2. $\boxed{\nabla_\eta^2 \log \mathcal{Z}(\eta) = F_\eta}$ (Fisher information matrix)

3.3. $\boxed{F_\eta = \nabla_\eta \xi = J_{\xi,\eta}}$ (Jacobian of mapping $\eta \to \xi$)

We will skip the details for now, as the Fisher information matrix will be covered in lecture 3. Also, discussion on these identities are in Chapter 3 of the course notes.

# Generalized Linear Models



Consider two familiar models

**Linear regression:** $\mathcal{L} = \frac{1}{2N} \sum_{i=1}^{N} (t_i - y_i)^2$, where $y_i = \boldsymbol{w}^\top \boldsymbol{x}_i$

$$\implies \nabla_{\boldsymbol{w}} \mathcal{L} = -\frac{1}{N} \sum_{i=1}^{N} (t_i - y_i) x_i$$

**Logistic regression:** $\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} -t_i \log y_i - (1 - t_i) \log(1 - y_i)$, where $y_i = \sigma(\boldsymbol{w}^\top \boldsymbol{x}_i)$

$$\implies \nabla_{\boldsymbol{w}} \mathcal{L} = -\frac{1}{N} \sum_{i=1}^{N} (t_i - y_i) x_i$$

Their gradients have the same form (!!!) Why?

# Generalized Linear Models

- Linear regression, logistic regression, softmax regression all belong to a broader class of models called **generalized linear models (GLM)**.
- GLM is derived from the exponential family.

# Generalized Linear Models

Consider the linear model with features:

$$z = \boldsymbol{w}^\top \phi(\boldsymbol{x})$$
$$y = a(z) \qquad \text{(activation)}$$

Assume the labels are distributed according to the exponential family (implied in the loss function)

$$p(t|\eta) = \frac{1}{\mathcal{Z}(\eta)} h(t) \exp(\eta^\top u(t))$$

We focus on a special case of the exponential family where $u(t) = t$.

$$p(t|\eta) = \frac{1}{\mathcal{Z}(\eta)} h(t) \exp(\eta^\top t)$$

# Generalized Linear Models

$$z = \boldsymbol{w}^\top \phi(\boldsymbol{x}), \quad y = a(z)$$

$$p(t|\eta) = \frac{1}{\mathcal{Z}(\eta)} h(t) \exp(\eta^\top t)$$

Recall, the moments can be computed by differentiating the partition function:

$$\nabla_\eta \log \mathcal{Z}(\eta) = \mathbb{E}_{t \sim p(t|\eta)}[u(t)] = \mathbb{E}_{t \sim p(t|\eta)}[t] = y \quad \text{(Prediction probability)}$$

There is a 1-to-1 mapping between $\eta \leftrightarrow y$. Let

$$\eta = \psi(y)$$

# Generalized Linear Models

$$z = \boldsymbol{w}^\top \phi(\boldsymbol{x}), \quad y = a(z)$$
$$\eta = \psi(y)$$
$$p(t|\eta) = \frac{1}{\mathcal{Z}(\eta)} h(t) \exp(\eta^\top t)$$

Gradient of log-likelihood w.r.t weights $\boldsymbol{w}$:

$$\frac{\partial}{\partial \boldsymbol{w}} \sum_{i=1}^{N} \log p(t_i|\eta_i) = \sum_{i=1}^{N} \frac{\partial}{\partial \eta_i} \log p(t_i|\eta_i) \frac{\partial \eta_i}{\partial y_i} \frac{\partial y_i}{\partial z_i} \frac{\partial z_i}{\partial \boldsymbol{w}}$$

$$= \sum_{i=1}^{N} (u(t_i) - \mathbb{E}_{t_i \sim p(t_i|\eta_i)}[u(t_i)]) \psi'(y_i) a'(z_i) \phi(\boldsymbol{x}_i)$$

$$= \sum_{i=1}^{N} (t_i - y_i) \psi'(y_i) a'(z_i) \phi(\boldsymbol{x}_i)$$

This greatly simplifies if we choose $a = \psi^{-1}$.

$$\eta = \psi(a(z)) = \psi(\psi^{-1}(z)) = z \implies \psi'(y)a'(z) = \frac{\partial \eta}{\partial y} \frac{\partial y}{\partial z} = 1$$

# Generalized Linear Models

To summarize, when the following conditions are met:

▶ Linear model with activation

$$z = \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}), \quad y = a(z)$$

▶ The label distribution belongs to the exponential family

$$p(t|\eta) = \frac{1}{\mathcal{Z}(\eta)} h(t) \exp(\eta^\top t)$$

where $y = \nabla_\eta \log \mathcal{Z}(\eta)$, and $\eta = \psi(y)$

▶ Activation is chosen as:

$$a(\cdot) = \psi^{-1}(\cdot)$$

Then we have:

$$\frac{\partial}{\partial \boldsymbol{w}} \sum_{i=1}^{N} \log p(t_i|\eta) = \sum_{i=1}^{N} (t_i - y_i) \boldsymbol{\phi}(\boldsymbol{x_i})$$

# GLM example: logistic regression

Cross-entropy loss (negative log-likelihood):

$$\mathcal{L} = -\log p(t|y) = -t \log y - (1-t) \log(1-y)$$

Corresponding label distribution: Bernoulli

$$p(t|y) = y^t (1-y)^{1-t}$$
$$= (1-y) \exp\{(\log \frac{y}{1-y})t\}$$

We have

$$\eta = \log \frac{y}{1-y} = \psi(y)$$

Then we should choose activation:

$$a(z) = \psi^{-1}(z) = \sigma(z) \quad \checkmark$$

# GLM example: linear regression

Squared loss (negative log-likelihood):

$$\mathcal{L} = -\log p(t|y) = \frac{1}{2}(t - y)^2$$

Corresponding label distribution: Gaussian (with fixed $\sigma$, WLOG assume $\sigma = 1$)

$$p(t|y) = \frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}(t - y)^2\}$$

We have

$$\eta = y = \psi(y)$$

Then we should choose activation:

$$a(z) = \psi^{-1}(z) = z \quad \checkmark$$

# Summary

Exponential family

$$p(\boldsymbol{x}|\boldsymbol{\eta}) = \frac{1}{\mathcal{Z}(\boldsymbol{\eta})} h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{u}(\boldsymbol{x}))$$

▶ Many common distributions belong to this family (Bernoulli, multinomial, Gaussian, Poisson, gamma, ...)
▶ Sufficient statistics for maximum-likelihood estimation
▶ Many useful identities stemming from $\mathcal{Z}(\boldsymbol{\eta})$
  ▶ Moments & empirical moments, MLE as moment matching
  ▶ Convenient way to compute the Fisher information matrix $\boldsymbol{F_\eta}$
▶ Used to derive the generalized linear models (GLM)