

# Scalable Generalized Linear Models

Check-in code: XX-XX-XX

Shuo Zhou, PhD

COM6012 Scalable Machine Learning  
05.03.2025



University of  
**Sheffield**

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

Generalized linear models in Spark ML

References and recommended reading

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

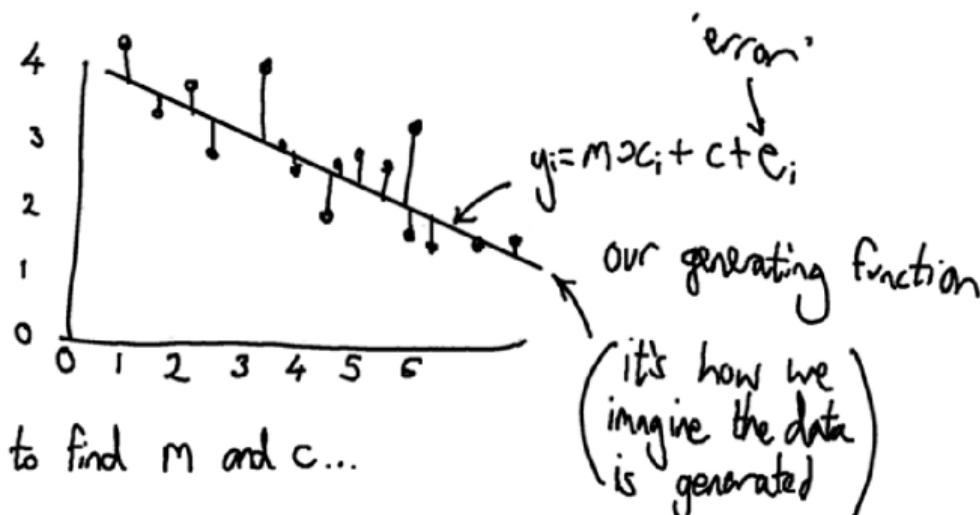
Generalized linear models in Spark ML

References and recommended reading

# Linear regression (I)

Recall COM4509/6509 (2024)

We have a model: Our data is from  
a linear, 1st-order polynomial with (Gaussian) noise  
added:



We need to find  $m$  and  $c$ ...

## Linear regression (II) - linear least squares

Recall COM4509/6509 (2024)

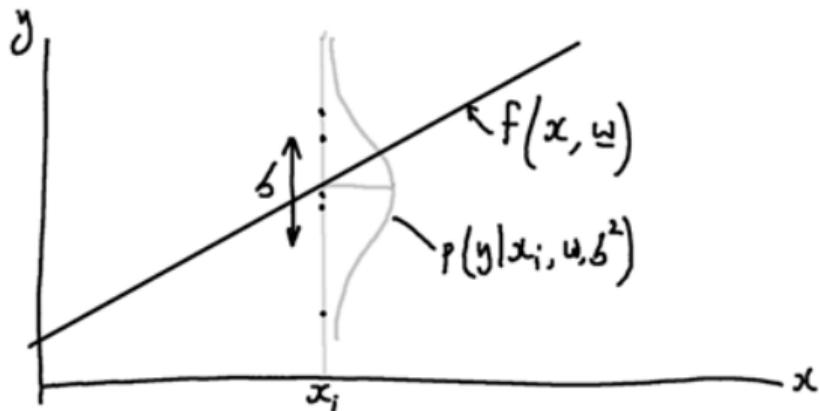
We need a way to say how good (or bad) a choice of  $m$  and  $c$  is.

We will use the sum squared error. Later we will see why this is a good choice.

$$E = \sum_{i=1}^N \left[ \underbrace{(mx_i + c)}_{\text{our prediction}} - \underbrace{y_i}_{\text{true label}} \right]^2$$

## Linear regression (III)

Recall COM4509/6509 (2024)



$$P(y_i|\mathbf{x}_i, \mathbf{w}, \sigma^2) = \mathcal{N}(y_i|\mathbf{w}^\top \mathbf{x}_i, \sigma^2)$$

$$\mathbb{E}(y_i|\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i, y \in \mathbb{R}$$

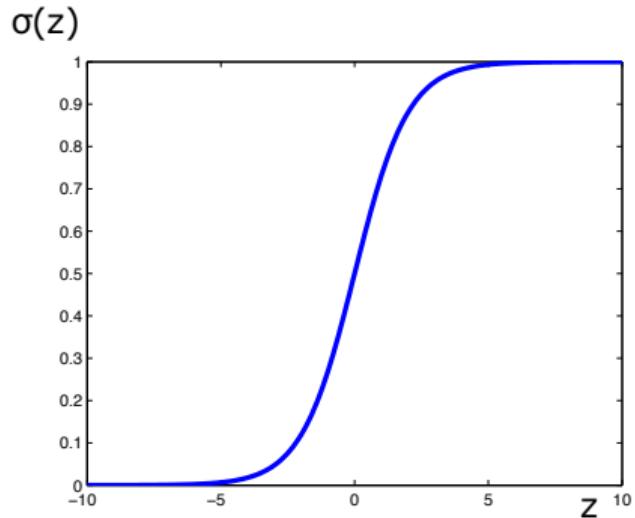
# Logistic regression

Recall  $p(y_i|\mathbf{x}_i, \mathbf{w}) = \text{Ber}(y_i|\sigma(\mathbf{w}^\top \mathbf{x}_i))$

$$P(Y=1) = \sigma(z) = \frac{1}{1 + \exp(-z)} = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

$$\max_{\mathbf{w}} \prod_{i=1}^N p(y_i|\mathbf{w}, \mathbf{x}_i) \rightarrow \min_{\mathbf{w}} \sum_{i=1}^N -\log p(y_i|\mathbf{w}, \mathbf{x}_i)$$

$$\mathbb{E}(y_i|\mathbf{x}_i) = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{\exp(\mathbf{w}^\top \mathbf{x}_i)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_i)}, y_i \in \{0, 1\}.$$



## Count data

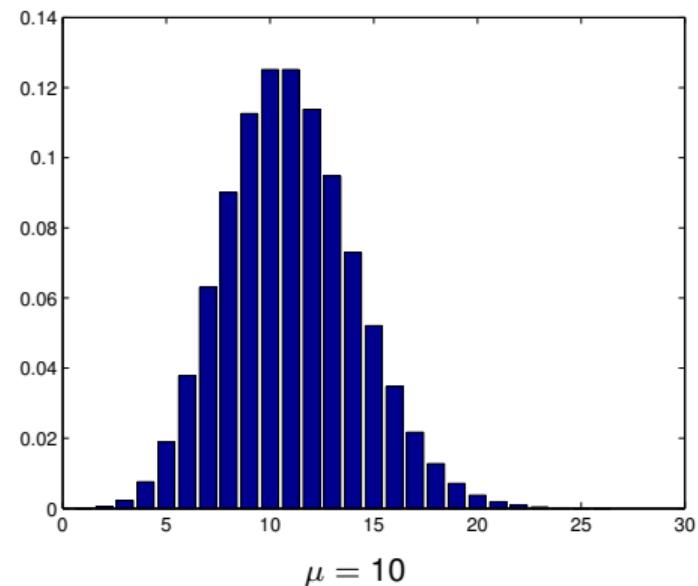
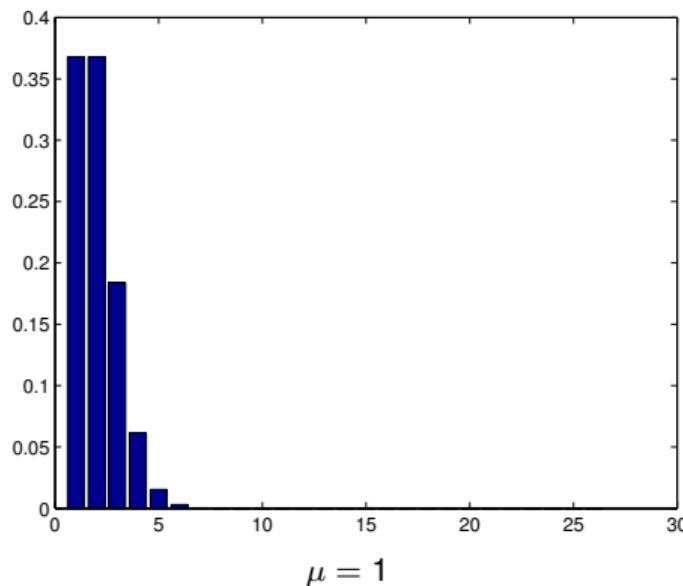
- Non-negative integer values, e.g. number of traffic incidences,  $Y \in \mathbb{N} = 0, 1, 2, 3, \dots, \infty$ .
  
- Discussion: Can we use linear regression or logistic regression for count data? Why?

## Poisson distribution

- A discrete probability distribution and the probability mass function (PMF) is

$$\text{Poi}(y|\mu) = \exp(-\mu) \frac{\mu^y}{y!},$$

where  $\mu > 0$  is called the *rate parameter*, and  $y \in \{0, 1, 2, \dots\}$ .



# Poisson regression

- Definition:

$$p(y_i | \mathbf{x}_i, \mathbf{w}) = \text{Poi}(y_i | \exp(\mathbf{w}^\top \mathbf{x}_i)),$$

where  $\text{Poi}(y|\mu) = \exp(-\mu) \frac{\mu^y}{y!}$

- The log probability density function (PDF):

$$\log p(y_i | \mathbf{x}_i, \mathbf{w}) = y_i \log \mu_i - \mu_i - \log(y_i!),$$

where  $\mu_i = \mathbb{E}(y_i | \mathbf{x}_i) = \exp(\mathbf{w}^\top \mathbf{x}_i)$

- Likelihood:

$$p(y_1, y_2, \dots, y_N | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N; \mathbf{w}) = \prod_{i=1}^N \exp(-\exp(\mathbf{w}^\top \mathbf{x}_i)) \frac{\exp(y_i \mathbf{w}^\top \mathbf{x}_i)}{y_i!}$$

# Generalized form?

- Linear regression:

$$\mathbb{E}(y|\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

- Logistic regression

$$\mathbb{E}(y|\mathbf{x}) = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- Poisson regression

$$\mathbb{E}(y|\mathbf{x}) = \exp(\mathbf{w}^\top \mathbf{x})$$

- General form?

$$g(\mathbb{E}(y|\mathbf{x})) = \mathbf{w}^\top \mathbf{x}, \text{ or } \mathbb{E}(y|\mathbf{x}) = g^{-1}(\mathbf{w}^\top \mathbf{x})$$

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

Generalized linear models in Spark ML

References and recommended reading

# Introduction

- By now you are familiar with different types of probability distributions: the Gaussian, the Bernoulli, the Poisson, etc.
- These are members of a broader class of distributions known as the exponential family.

## Why the exponential family is important?

- It can be shown that the exponential family is the only family of distributions with finite-sized sufficient statistics.
- The exponential family is the only family of distributions for which conjugate priors exist.
- The exponential family can be shown to be the family of distributions that makes the least set of assumptions subject to some user-chosen constraints.
- The exponential family is at the core of generalized linear models.

## Definition

- It is said that a probability density function (PDF) or a probability mass function (PMF)  $p(\mathbf{y}|\boldsymbol{\theta})$ , with  $\mathbf{y} \in \mathbb{R}^p$  and  $\boldsymbol{\theta} \in \mathbb{R}^d$ , is in the **exponential family** if it is of the form

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{y}) \exp [\boldsymbol{\theta}^\top \mathcal{T}(\mathbf{y})],$$

where

$$Z(\boldsymbol{\theta}) = \int_{\mathbf{y}} h(\mathbf{y}) \exp [\boldsymbol{\theta}^\top \mathcal{T}(\mathbf{y})] d\mathbf{y}.$$

- $h(\mathbf{y})$  is a scaling constant, often 1.
- $\boldsymbol{\theta}$  are known as the **natural parameters** or **canonical parameters**.
- $\mathcal{T}(\mathbf{y}) \in \mathbb{R}^d$  is called a vector of **sufficient statistics**.
- $Z(\boldsymbol{\theta})$  is known as the **partition function**.

## Definition (II)

- Distributions in the exponential family can also be expressed as

$$p(\mathbf{y}|\boldsymbol{\theta}) = h(\mathbf{y}) \exp [\boldsymbol{\theta}^\top \mathcal{T}(\mathbf{y}) - A(\boldsymbol{\theta})],$$

where

$$A(\boldsymbol{\theta}) = \log Z(\boldsymbol{\theta}).$$

- $A(\boldsymbol{\theta})$  is called the **log partition function** or **cumulant function**.
- If  $\mathcal{T}(\mathbf{y}) = \mathbf{y}$ , we say it is a **natural exponential family**.

## Example: Bernoulli

- For the Bernoulli distribution,  $y \in \{0, 1\}$ , and we have

$$\begin{aligned}\text{Ber}(y|\mu) &= \mu^y(1-\mu)^{1-y} \\ &= \exp[y \log \mu + (1-y) \log(1-\mu)] \\ &= \exp\left[y \log\left(\frac{\mu}{1-\mu}\right) + \log(1-\mu)\right], \\ &= h(y) \exp[\boldsymbol{\theta}^\top \mathcal{T}(y) - A(\boldsymbol{\theta})],\end{aligned}$$

where

- $\boldsymbol{\theta} = \log\left(\frac{\mu}{1-\mu}\right)$ , known as **log-odds ratio**
- $\mathcal{T}(y) = y$
- $A(\boldsymbol{\theta}) = -\log(1-\mu)$
- $h(y) = 1$

## Example: univariate Gaussian distribution

- The univariate Gaussian can be written in an exponential family form as

$$\begin{aligned}\mathcal{N}(y|\mu, \sigma^2) &= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left[-\frac{1}{2\sigma^2}(y-\mu)^2\right] \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left[-\frac{1}{2\sigma^2}y^2 + \frac{\mu}{\sigma^2}y - \frac{1}{2\sigma^2}\mu^2\right] \\ &= \frac{1}{Z(\theta)} \exp(\theta^\top \mathcal{T}(y)),\end{aligned}$$

where

- $\theta = \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix}, \mathcal{T}(y) = \begin{bmatrix} y \\ y^2 \end{bmatrix}$

- $Z(\theta) = \sqrt{2\pi}\sigma \exp\left\{\frac{\mu^2}{2\sigma^2}\right\}$

- $h(y) = 1$

## Example: Poisson distribution

- As a member of the exponential family,  $\text{Poi}(y|\mu) = \exp(-\mu) \frac{\mu^y}{y!}$  can be written as

$$\text{Poi}(y|\mu) = \frac{h(y)}{Z(\theta)} \exp(\theta y),$$

where  $\theta = \log \mu$ ,  $h(y) = \frac{1}{y!}$ , and  $Z(\theta) = \exp(\mu)$ ,  $A(\theta) = \mu$ .

- Recall that the expected value of  $y$  is equal to  $\mathbb{E}[y] = \mu$ .
- Then, the mean parameter  $\mu$  can be recovered from the canonical parameter using

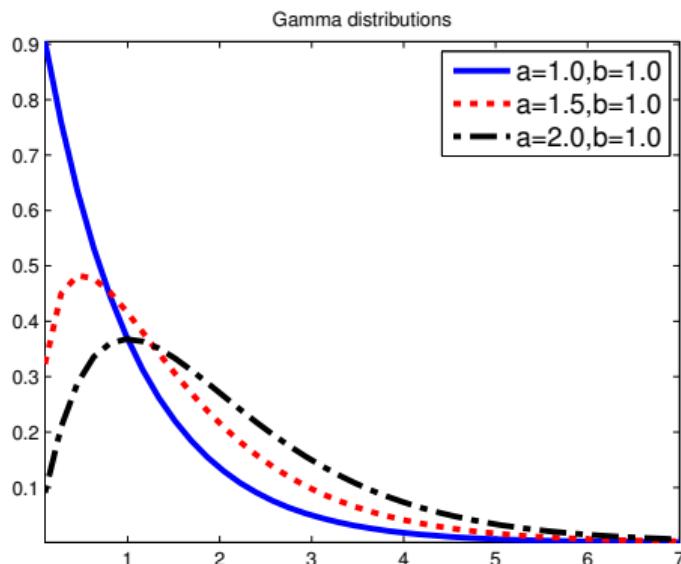
$$\mu = \exp(\theta).$$

## Example: Gamma distribution (I)

- The Gamma distribution follows as

$$Ga(y|a, b) = \frac{b^a}{\Gamma(a)} y^{a-1} \exp(-by),$$

where  $a > 0$  (shape), and  $b > 0$  (rate).  $\Gamma(a) = \int_0^\infty u^{a-1} e^{-u} du$  is the Gamma function.



## Example: Gamma distribution (II)

- As a member of the exponential family, it can be written as

$$\text{Ga}(y|a, b) = \frac{1}{Z(\theta)} \exp(\theta^\top \mathcal{T}(y)),$$

where

- $\theta = \begin{bmatrix} a-1 \\ -b \end{bmatrix}$ ,  $\mathcal{T}(y) = \begin{bmatrix} \log y \\ y \end{bmatrix}$
- $Z(\theta) = \frac{\Gamma(a)}{b^a}$ ,  $A(\theta) = \log \Gamma(a) - a \log b$
- $h(y) = 1$

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

Generalized linear models in Spark ML

References and recommended reading

## Definition

- Linear, logistic, and Poisson regression are examples of generalized linear models (GLM).
- These are models in which the output is in the exponential family.
- The mean parameters are a linear combination of the inputs, passed through a possibly nonlinear function, such as the logistic function.

## General form (I)

- We want to model the relationship between a response variable  $y_i$ , and an input vector  $\mathbf{x}_i$ .
- Let us first consider the case of an unconditional distribution for the response variable

$$p(y_i|\theta, \sigma^2) = \exp \left[ \frac{y_i\theta - A(\theta)}{\sigma^2} + c(y_i, \sigma^2) \right],$$

where  $\theta$  is the natural parameter,  $A(\cdot)$  is the log-partition function,  $\sigma^2$  is the **dispersion parameter**, usually,  $\sigma^2 = 1$ , and  $c$  is the normalization constant.

- The expression for  $p(y_i|\theta, \sigma^2)$  looks similar to the exponential family.

## General form (II)

- For example, in logistic regression,  $\theta$  is the log-odds ratio

$$\theta = \log \left( \frac{\mu}{1 - \mu} \right),$$

where  $\mu = \mathbb{E}[y] = P(y = 1)$  is the mean parameter.

- To convert from the mean parameter  $\mu$  to the natural parameter  $\theta$ , we can use a function  $\psi$ ,  $\theta = \psi(\mu)$ .
- $\psi$  is uniquely determined by the form of the exponential family distribution.
- The mapping is invertible, so that  $\mu = \psi^{-1}(\theta)$ .

## Link function (I)

- Now let us add inputs or covariates.
- We first define a linear function of the inputs  $\eta_i = \mathbf{w}^\top \mathbf{x}_i$ .
- We now make the mean of the distribution be some invertible monotonic function of this linear combination.
- By convention, this function, known as the **mean function**, is denoted by  $g^{-1}(\cdot)$ , so

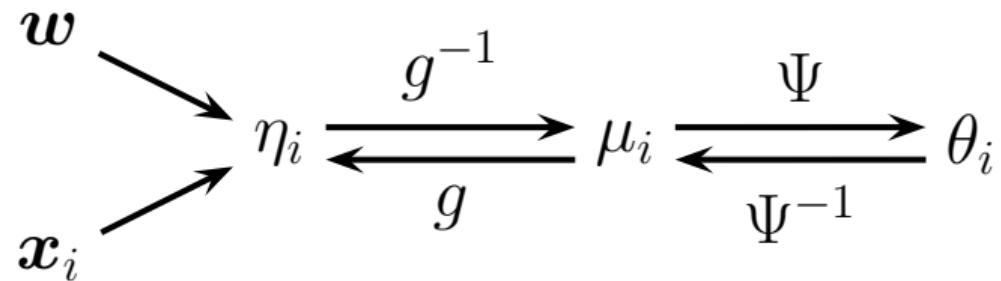
$$\mathbb{E}(y_i | \mathbf{x}_i) = \mu_i = g^{-1}(\eta_i) = g^{-1}(\mathbf{w}^\top \mathbf{x}_i).$$

- The inverse of the mean function, namely  $g(\cdot)$ , is called the **link function**.

## Link function (II)

- We are free to choose almost any function we like for  $g$ , so long as it is invertible, and so long as  $g^{-1}$  has the appropriate range.
- For example, in logistic regression, we set  $\mu_i = g^{-1}(\eta_i) = \sigma(\eta_i)$ .

## Relationships between functions



$g^{-1}()$  is the **mean function**.  $g()$  is the **link function**.

## GLM with canonical link function

- One particularly simple form of link function is to use  $g = \psi$ .
- This is called the **canonical link function**.
- In this case  $\theta_i = \eta_i = \mathbf{w}^\top \mathbf{x}_i$ , so the model becomes

$$p(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2) = \exp \left[ \frac{y_i \mathbf{w}^\top \mathbf{x}_i - A(\mathbf{w}^\top \mathbf{x}_i)}{\sigma^2} + c(y_i, \sigma^2) \right].$$

## Canonical link functions $g = \psi$ for some GLMs

Distribution	Link $g(\mu)$	$\theta = \psi(\mu)$	$\mu = \psi^{-1}(\theta)$
$\mathcal{N}(\mu, \sigma^2)$	identity	$\theta = \mu$	$\mu = \theta$
Ber( $\mu$ )	logit	$\theta = \log\left(\frac{\mu}{1-\mu}\right)$	$\mu = \sigma(\theta)$
Poi( $\mu$ )	log	$\theta = \log(\mu)$	$\mu = \exp(\theta)$
Ga( $a, b$ )	inverse	$\theta = \mu^{-1}$	$\mu = \theta^{-1}$ .

## Mean and variance of the response variable

- Recall  $A(\cdot)$  as the *cumulant function* (or log partition function)
- It can be shown that

$$\begin{aligned}\mathbb{E}[y|\mathbf{x}_i, \mathbf{w}, \sigma^2] &= \mu_i = A'(\theta_i) \\ \text{var}[y|\mathbf{x}_i, \mathbf{w}, \sigma^2] &= \sigma_i^2 = A''(\theta_i)\sigma^2.\end{aligned}$$

- Note the dispersion parameter  $\sigma^2$  is not the variance of  $y$ .

## Example: linear regression

- In linear regression, the response variable follows a normal distribution,

$$\begin{aligned} p(y_i|\mu_i, \sigma^2) &= \mathcal{N}(y_i|\mu_i, \sigma^2) \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left[-\frac{1}{2\sigma^2}(y_i - \mu_i)^2\right] \\ &= \exp\left[\frac{y_i\mu_i - \frac{\mu_i^2}{2}}{\sigma^2} - \frac{y_i^2}{2\sigma^2} + \log\left(\frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}}\right)\right] \\ &= \exp\left[\frac{y_i\mu_i - \frac{\mu_i^2}{2}}{\sigma^2} - \frac{1}{2}\left(\frac{y_i^2}{\sigma^2} + \log(2\pi\sigma^2)\right)\right]. \end{aligned}$$

- For linear regression,  $y_i \in \mathbb{R}$ . The link function is the identity  $\theta_i = \mu_i = \mathbf{w}^\top \mathbf{x}_i$ .
- Recall GLM with canonical link functions:  $p(y_i|\mathbf{x}_i, \mathbf{w}, \sigma^2) = \exp\left[\frac{y_i\mathbf{w}^\top \mathbf{x}_i - A(\mathbf{w}^\top \mathbf{x}_i)}{\sigma^2} + c(y_i, \sigma^2)\right]$ .
- With  $A(\mu_i) = \mu_i^2/2$ ,  $\mathbb{E}[y_i] = \mu_i$ , and  $\text{var}[y_i] = \sigma^2$ .

## Example: logistic regression

- In logistic regression, the response variable follows a Bernoulli distribution

$$\begin{aligned} p(y_i|\mu_i, \sigma^2) &= \mu_i^{y_i}(1 - \mu_i)^{1-y_i} \\ &= \exp \left[ \log \left( \frac{\mu_i}{1 - \mu_i} \right) y_i - (-\log(1 - \mu_i)) \right]. \end{aligned}$$

- The link function is the logit,  $g(\mu) = \theta = \log \left( \frac{\mu_i}{1 - \mu_i} \right) = \mathbf{w}^\top \mathbf{x}_i$ .
- Recall  $p(y_i|\mathbf{x}_i, \mathbf{w}, \sigma^2) = \exp \left[ \frac{y_i \mathbf{w}^\top \mathbf{x}_i - A(\mathbf{w}^\top \mathbf{x}_i)}{\sigma^2} + c(y_i, \sigma^2) \right]$ .
- With  $A(\theta_i) = -\log(1 - \sigma(\theta_i))$ ,  $\mathbb{E}[y_i] = \sigma(\theta_i)$ , and  $\text{var}[y_i] = \sigma(\theta_i)(1 - \sigma(\theta_i))$ .

## Example: Poisson regression

- In Poisson regression, the response variable follows a Poisson distribution

$$p(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2) = \exp[y_i \log(\mu_i) - \mu_i - \log(y_i!)].$$

- The link function is  $\log$ ,  $\log \mu_i = \mathbf{w}^\top \mathbf{x}_i$ .

- Recall  $p(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2) = \exp\left[\frac{y_i \mathbf{w}^\top \mathbf{x}_i - A(\mathbf{w}^\top \mathbf{x}_i)}{\sigma^2} + c(y_i, \sigma^2)\right]$ .

- With  $A(\theta_i) = \exp(\theta_i)$ ,  $\mathbb{E}[y_i] = \exp(\theta_i) = \mu_i$ .

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

Generalized linear models in Spark ML

References and recommended reading

## Log-likelihood for a GLM

- ❑ One of the appealing properties of GLMs is that they can be fit using the same methods that we used to fit logistic regression.
- ❑ In particular, the log-likelihood has the following form

$$\mathcal{L}(\mathbf{w}) = \frac{1}{\sigma^2} \sum_{i=1}^N \mathcal{L}_i = \frac{1}{\sigma^2} \sum_{i=1}^N [\theta_i y_i - A(\theta_i)].$$

where  $\mathcal{L}_i = \theta_i y_i - A(\theta_i)$ .

## Gradient for the log-likelihood

- We can compute the gradient vector using the chain rule as follows

$$\begin{aligned}\frac{\partial \mathcal{L}_i}{\partial w_j} &= \frac{d\mathcal{L}_i}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial w_j} \\ &= (y_i - A'(\theta_i)) \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} x_{i,j} \\ &= (y_i - \mu_i) \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} x_{i,j}.\end{aligned}$$

- If we use a canonical link,  $\theta_i = \eta_i$ , this simplifies to

$$\mathbf{g}(\mathbf{w}) = \frac{1}{\sigma^2} \left[ \sum_{i=1}^N (y_i - \mu_i) \mathbf{x}_i \right] = \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}),$$

where  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_N]^\top$ .

- This can be used inside a (stochastic) gradient descent procedure.

## Hessian for the log-likelihood

- For improved efficiency, we could use a second-order method.
- If we use a canonical link, the Hessian is given by

$$\mathbf{H}(\mathbf{w}) = -\frac{1}{\sigma^2} \sum_{i=1}^N \frac{d\mu_i}{d\theta_i} \mathbf{x}_i \mathbf{x}_i^\top = -\frac{1}{\sigma^2} \mathbf{X}^\top \Sigma \mathbf{X},$$

where  $\Sigma = \text{diag} \left( \frac{d\mu_1}{d\theta_1}, \dots, \frac{d\mu_N}{d\theta_N} \right)$ .

- This can be used inside the Iterative Reweighted Least Squares (IRLS) algorithm.

## Iterative reweighted least squares (IRLS) algorithm

- ❑ The Iterative Reweighted Least Squares algorithm is a particular case of the Newton's method.
- ❑ The updated parameters are obtained by iteratively solving a weighted least squares problem.

## A least squares problem

- Recall that a least squares (*LS*) problem refers to

$$LS(\mathbf{w}) = \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2,$$

for a dataset  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N = \{\mathbf{X}, \mathbf{y}\}$ .

- It can be shown that the vector  $\mathbf{w}$  that minimizes  $LS(\mathbf{w})$  is given as

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

## A weighted least squares problem

- A weighted least squares (*WLS*) problem refers to

$$WLS(\mathbf{w}) = \min_{\mathbf{w}} \sum_{i=1}^N r_i(y_i - \mathbf{w}^\top \mathbf{x}_i)^2,$$

for a dataset  $\mathcal{D} = \{\mathbf{x}_i, y_i, r_i\}_{i=1}^N = \{\mathbf{X}, \mathbf{R}, \mathbf{y}\}$ , with  $\mathbf{R} = \text{diag}(r_1, \dots, r_N)$ .

- It can be shown that the vector  $\mathbf{w}$  that minimizes  $WLS(\mathbf{w})$  is given as

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{R} \mathbf{y}.$$

## Iterative reweighted least squares problem

- Newton's method for the log-likelihood of the GLM follows as

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}_k^{-1} \mathbf{g}_k \\ &= \mathbf{w}_k + (\mathbf{X}^\top \Sigma_k \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}_k) \\ &= (\mathbf{X}^\top \Sigma_k \mathbf{X})^{-1} [\mathbf{X}^\top \Sigma_k \mathbf{X} \mathbf{w}_k + \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}_k)] \\ &= (\mathbf{X}^\top \Sigma_k \mathbf{X})^{-1} \mathbf{X}^\top \Sigma_k \mathbf{z}_k,\end{aligned}$$

where  $\mathbf{z}_k = \mathbf{X} \mathbf{w}_k + \Sigma_k^{-1} (\mathbf{y} - \boldsymbol{\mu}_k)$  is known as the **working response**.

- Recall  $WLS(\mathbf{w})$ :  $\mathbf{w} = (\mathbf{X}^\top \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{R} \mathbf{y}$ . At iteration  $k$ , the solution for  $\mathbf{w}_{k+1}$  has a similar form to the solution for a weighted least squared problem replacing  $\mathbf{R}$  for  $\Sigma_k$ , and  $\mathbf{y}$  for  $\mathbf{z}_k$ .
- The name IRLS is due to at each iteration, we solve a weighted least squares problem, where the weight matrix  $\Sigma_k$  changes at each iteration.

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

Generalized linear models in Spark ML

References and recommended reading

## GeneralizedLinearRegression()

- ❑ It uses IRLS (Iterative Reweighted Least Squares) for optimization.
- ❑ It only allows  $\ell_2$  regularization.
- ❑ Spark currently only supports up to **4096 features** through its GeneralizedLinearRegression interface.
- ❑ It will throw an exception if this constraint is exceeded.

## GLM available in Spark

- It includes the following families

Family	Response type	Supported links
Gaussian	Continuous	Identity*, Log, Inverse
Binomial	Binary	Logit*, Probit, CLogLog
Poisson	Count	Log*, Identity, Sqrt
Gamma	Continuous	Inverse*, Identity, Log
Tweedie	Zero-inflated continuous	Power link function

where \* stands for canonical link.

- For any random variable  $Z$  that obeys a Tweedie distribution, the variance  $\text{var}(Z)$  relates to the mean  $\mathbb{E}(Z)$  by the power law:  $\text{var}(Z) = a\mathbb{E}(Z)^p$ , where  $a$  and  $p$  are positive constants.
- The parameters are set using `family` and `link`.

## Parameters to adjust

- ❑ **maxIter**: max number of iterations.
- ❑ **regParam**: regularization parameter ( $\geq 0$ ).
- ❑ **family**: name of the family which describes the label distribution to be used in the model.
- ❑ **link**: name of link function which provides the relationship between the linear predictor and the mean of the distribution function.

## Particular cases: `LinearRegression()`

- ❑ If your ‘family’ is Gaussian and the link function is the ‘identity’, your model is just equivalent to linear regression.
- ❑ My recommendation is that you use `LinearRegression()` instead of `GeneralizedLinearRegression()`.
- ❑ `LinearRegression()` allows for  $\ell_1$ ,  $\ell_2$  and elastic net regularization through L-BFGS or OWL-QN.

## Particular cases: LogisticRegression()

- ❑ If your ‘family’ is Binomial and the link function is ‘logit’, your model is equivalent to logistic regression.
- ❑ My recommendation is that you use `LogisticRegression()` instead of `GeneralizedLinearRegression()`.
- ❑ `LogisticRegression()` allows for  $\ell_1$ ,  $\ell_2$  and elastic net regularization through L-BFGS or OWL-QN.

# Contents (Check-in code: XX-XX-XX)

Linear regression, logistic regression, and Poisson regression

The exponential family

Generalized linear models (GLM)

Iteratively reweighted least squares (IRLS)

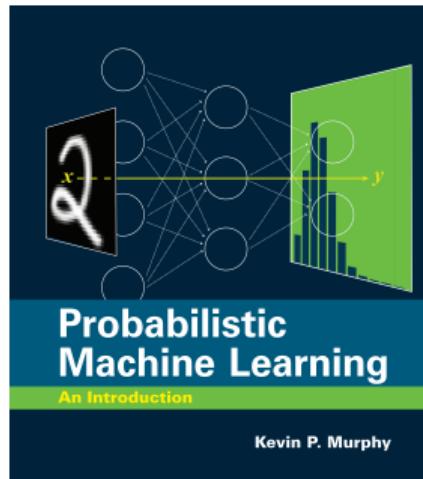
Generalized linear models in Spark ML

References and recommended reading

# References and recommended reading

Book: *Probabilistic Machine Learning: An Introduction* by Kevin P Murphy, 2022.

- Section 3.4 The exponential family, pp. 90-93 (book) / pp. 93-96 (draft .pdf file).
- Chapter 12 Generalized linear models, pp. 409-415 (book) / pp. pp 415-421 (draft .pdf file).



# References and recommended reading

Youtube video: *Generalized Linear Models in Spark MLlib and SparkR* by Xiangrui Meng (from Databricks).

The screenshot shows a YouTube video player interface. At the top, there is a navigation bar with a menu icon, the YouTube logo, and a search bar. Below the navigation bar is the video thumbnail, which features a green landscape background and the title 'Generalized Linear Models in Spark MLlib and SparkR' in white text. Below the title, it says 'Xiangrui Meng' and 'joint with Joseph Bradley, Eric Liang, Yanbo Liang (MiningLamp), DB Tsai (Netflix), et al.'. It also includes the date '2016/02/17 - Spark Summit East' and the 'databricks' logo. To the right of the thumbnail, there is a vertical banner for 'SPARK SUMMIT EAST' with the text 'DATA SCIENCE AND ENGINEERING AT SCALE' and 'FEBRUARY 16-18, 2016 NEW YORK CITY'. The main video player area shows a person speaking at a podium. The video progress bar indicates it is at 0:21/29:05. Below the video player, the video title 'Generalized Linear Models in Spark MLlib and SparkR' is displayed, along with '2,521 views', '11 likes', '1 dislike', 'SHARE', 'SAVE', and a 'More' options menu. The bottom of the screen shows standard YouTube navigation icons.

[https://www.youtube.com/watch?v=PSZW6hcQ\\_7w](https://www.youtube.com/watch?v=PSZW6hcQ_7w)

## References and recommended reading

Website: [Spark Python API Documentation](#).

<https://spark.apache.org/docs/3.5.4/api/python/reference/index.html>

## References and recommended reading

Paper: [Map-Reduce for Machine Learning on Multicore](#) by C-T Chu et al. (2006).

Map-Reduce for Machine Learning on Multicore

**Cheng-Tao Chu \***      **Sang Kyun Kim \***      **Yi-An Lin \***  
chentao@stanford.edu    skkim38@stanford.edu    ianl@stanford.edu

**Yuan Yuan Yu \***      **Gary Bradski †**      **Andrew Y. Ng \***  
yuanyuan@stanford.edu    garybradski@gmail.com    ang@cs.stanford.edu

**Kunle Olukotun \***  
kunle@cs.stanford.edu

\* CS. Department, Stanford University 353 Serra Mall  
Stanford University, Stanford CA 94305-9025.

Rexee Inc.

### Abstract

We are at the beginning of the multicore era. Computers will have increasingly many cores (processors), but there is still no good programming framework for these architectures, and thus no simple and unified way for machine learning to take advantage of the potential speed up. In this paper, we develop a broadly applicable parallel programming method, one that is easily applied to *many* different learning algorithms. Our work is in distinct contrast to the tradition in machine learning of designing (often ingenious) ways to speed up a *single* algorithm at a time. Specifically, we show that algorithms that fit the Statistical Query model [15] can be written in a certain “summation form,” which allows them to be easily parallelized on multicore computers. We adapt Google’s map-reduce [7] paradigm to demonstrate this parallel speed up technique on a variety of learning algorithms, including locally weighted linear regression (LWLR), k-means, logistic regression (LR), naive Bayes (NB), SVM, ICA, PCA, gaussian discriminant analysis (GDA), EM, and backpropagation (NN). Our experimental results show basically linear speed up with an increasing number of processors.

## Acknowledgment

- ❑ Thanks to [Dr. Robert Loftin](#) for his contribution to this session in 2024.
  
- ❑ Thanks to [Dr. Mauricio Alvarez](#) for his contribution to this session from 2017 to 2022.