

Part 3:

Generalized linear models

Logistic regression

Approximate inference: Laplace method

---

### 3.1. Generalized linear models

$\vec{x} \in \mathbb{R}^d$  - input,  $y$  - target

$D = \{(\vec{x}_i, y_i)\}_{i=1}^n$  - a sample

Goal: get  $p(y | \vec{x})$

General model:

$\lambda = \lambda(\vec{x}; \vec{\Theta}) = \vec{x}^T \vec{\Theta}$  - linear link f-on

$$p(y | \vec{x}) = p(y | \lambda)$$

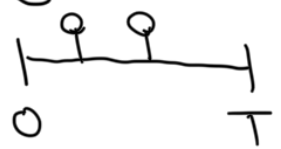
A. Linear regression

$$y \in \mathbb{R}, \quad \lambda = \vec{\Theta}^T \vec{x}$$

$y \sim N(\lambda, \sigma^2)$ ,  $\sigma^2$  is known / not important

B. Poisson regression

$y \in \mathbb{N} \cup \{0\}$  - number of events

 Intensity  $\exp(\lambda) > 0$

$$y \sim \text{Pois}(\exp(\lambda))$$

C. Logistic regression

$$y \in \{0, 1\}$$

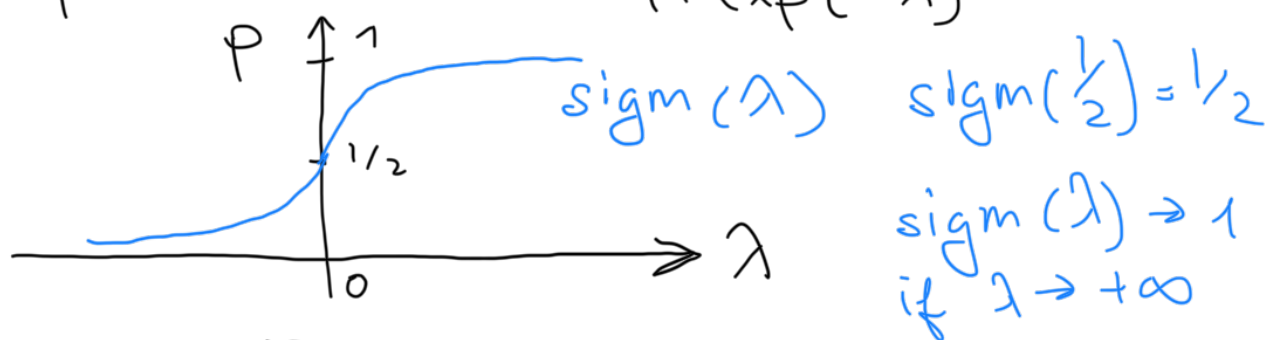
$y \sim \text{Bern}(p)$  - Bernoulli distrib.

1 1 1 0 1 ~

$y = 0$  with probability  $p$

$y = 1$  with probability  $1 - p$

$$p = \text{sigm}(\lambda) = \frac{1}{1 + \exp(-\lambda)} \in [0, 1]$$



$$\lambda = \vec{\Theta}^T \vec{x}$$

Data generation process:

$$\boxed{\vec{x}} \xrightarrow{\vec{\Theta}} \lambda \xrightarrow{\text{sigm}(\lambda)} p \rightarrow y \in \{0, 1\}$$

GLMs are:

- general
- powerful

How to estimate  $\vec{\Theta}$  from  $D$  for GLMs?

---

3.2. MLE - maximum likelihood estimate for logistic regression

$$y \in \{0, 1\}$$

$$p = p(y=0 | \vec{x}) = \mathcal{O}(\vec{x}^T \vec{\Theta}), \quad p_i = \mathcal{O}(\vec{x}_i^T \vec{\Theta})$$

$$p(y | \vec{x}, \vec{\Theta}) = p^y (1-p)^{1-y}$$

$$p(\vec{y} | X, \vec{\Theta}) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i},$$

$$\log p(\vec{y} | X, \vec{\Theta}) = \sum_{i=1}^n y_i \log p_i + (1-y_i) \log(1-p_i)$$

$$\frac{\partial}{\partial \theta} \log p(y|X, \vec{\theta}) = \sum_{i=1}^n (\underbrace{p_i - y_i}_{\substack{\text{depends on } \vec{\theta} \text{ in} \\ \text{nonlinear way}}}) \vec{x}_i = 0$$

### 3.3. Newton method for optimization

II-order method

$\vec{\theta}^{(0)}$  - initial guess

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} + \underbrace{H^{-1}}_{\substack{\text{instead of 'learning rate'} \\ \text{we use inverse Hessian matrix}}} \nabla \ln p(\vec{y}|X, \vec{\theta}) \quad (1)$$

instead of 'learning rate'

we use inverse Hessian matrix

$$H = \left\{ \frac{\partial^2 \ln p(\vec{y}|X, \vec{\theta})}{\partial \theta_i \partial \theta_j} \right\}$$

Statement We solve quadratic problem in one iteration via Newton's method

$$\blacksquare L(\vec{\theta}) = (\vec{\theta} - \vec{\mu})^T A (\vec{\theta} - \vec{\mu}) \rightarrow \min_{\vec{\theta}}, \quad A \succeq 0$$

$\vec{\theta}^{(0)}$  - initial guess

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - H^{-1} \left. \frac{\partial L(\vec{\theta})}{\partial \vec{\theta}} \right|_{\vec{\theta} = \vec{\theta}^{(0)}}$$

Now we need  $\frac{\partial L(\vec{\theta})}{\partial \vec{\theta}}$  and  $H$

$$\left\{ \begin{array}{l} \frac{\partial L(\vec{\theta})}{\partial \vec{\theta}} = 2A(\vec{\theta} - \vec{\mu}) \\ H = 2A, \quad H^{-1} = \frac{1}{2} A^{-1} \end{array} \right.$$

$$\begin{aligned} \vec{\theta}^{(1)} &= \vec{\theta}^{(0)} - \frac{1}{2} A^{-1} 2A(\vec{\theta}^{(0)} - \vec{\mu}) = \\ &= \vec{\theta}^{(0)} - \vec{\theta}^{(0)} + \vec{\mu} = \vec{\mu} \end{aligned}$$

Let's apply Newton's method for the optimization for parameters of Logistic regression

$$\ln p(\vec{y} | \vec{\Theta}) = \sum_{i=1}^n y_i \log p_i + (1-y_i) \log(1-p_i)$$

$$p_i = \frac{1}{1 + \exp(-\vec{x}_i^T \vec{\Theta})}; \quad \vec{x}_i^T \vec{\Theta} = t_i$$

$$\vec{p} = (p_1, \dots, p_n)^T$$

$$\begin{aligned} \frac{\partial p_i}{\partial \Theta_j} &= -\sigma^2(t_i) \exp(-t_i) (-x_{ij}) \\ &= \sigma(t_i) (1 - \sigma(t_i)) x_{ij} \\ &= p_i (1 - p_i) x_{ij} \end{aligned}$$

$$\begin{aligned} y_i = 1: \frac{\partial}{\partial \Theta_j} \ln p(y_i | \vec{\Theta}) &= \frac{1}{p_i} \cdot \frac{\partial p_i}{\partial \Theta_j} = \\ &= (1 - p_i) x_{ij} \end{aligned}$$

$$\begin{aligned} y_i = 0: \frac{\partial}{\partial \Theta_j} \ln p(y_i | \vec{\Theta}) &= -\frac{1}{1 - p_i} \frac{\partial p_i}{\partial \Theta_j} = \\ &= (0 - p_i) x_{ij} \end{aligned}$$

$$\nabla \ln p(\vec{y} | \vec{\Theta}) = (\vec{p} - \vec{y})^T X,$$

Hessian

$$H = X^T R X,$$

$$R = \text{diag}(\vec{p}(1-\vec{p}))$$

$$\vec{\Theta}^{(new)} = (X^T R X)^{-1} X^T R \vec{z}$$

$$R = \begin{pmatrix} p_1(1-p_1) & & 0 \\ & \ddots & \\ 0 & & p_n(1-p_n) \end{pmatrix}$$

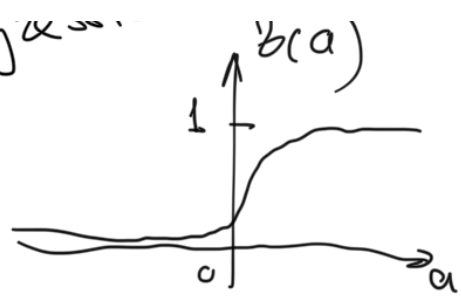
$$\vec{z} = X \vec{\Theta}^{(old)} - R^{-1}(\vec{p} - \vec{y})$$

$$p_i \in (0, 1) \\ p_i = \sigma(\vec{x}_i^T \vec{\Theta})$$

$$\vec{\Theta}^{(new)} = (X^T X)^{-1} X^T \vec{u} \quad \text{for linear regression}$$

$\cup_{i \in E} ( \quad ) \quad \cup \quad \text{reg} \alpha \dots$

$$\forall \vec{v}: \vec{v}^T H \vec{v} = \underbrace{\vec{v}^T X^T}_{\vec{v}'^T} R X \vec{v} = \underbrace{\vec{v}'^T}_{\vec{v}'^T} \underbrace{R}_{\vec{v}'} \vec{v}' = \sum_{i=1}^d z_{ii} (v_i')^2 \geq 0,$$



$\Rightarrow H$  is nonnegative-definite

$\Rightarrow$  the optimization problem is convex

### 3.4. Probit regression

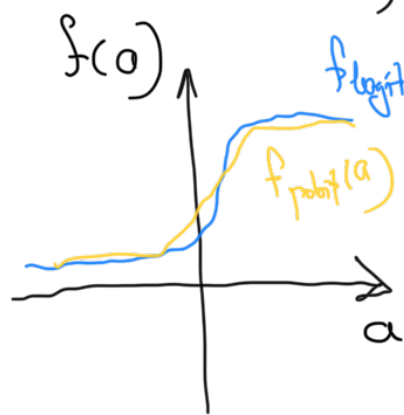
$$p(y=0 | a) = f(a), \quad a = \vec{\theta}^T \vec{x}$$

$$\begin{array}{l} a_i = \vec{\theta}^T \vec{x}_i \\ t_i \sim p(\cdot) \end{array} \rightarrow y_i = \mathbb{I}[a_i > t_i]$$

$$f(a) = \int_{-\infty}^a p(t) dt \quad (-\infty, a_i]$$

$$f(a) = \begin{cases} \Phi(a) = \int_{-\infty}^a \underbrace{N(t|0,1)}_{\text{no analytical formula for } \Phi(a)} dt & \text{probit regression} \\ p_{\text{log}}(t) = \frac{\exp(-t)}{(1+\exp(-t))^2} & \text{logistic regression} \end{cases}$$

$f(a) = \delta(a) = \frac{1}{1+\exp(-a)}$



### 3.5. Approximate Bayesian Inference: idea

$$p(\theta | \mathcal{D}) \propto p(\mathcal{D} | \theta) \pi(\theta)$$

$$p(\mathcal{D}) = ? = \int p(\mathcal{D} | \theta) \pi(\theta) d\theta$$

$$p(y | \mathcal{D}) = \int p(y | \vec{x}, \vec{\theta}) p(\vec{\theta} | \mathcal{D}) d\vec{\theta} = ?$$

$$p(\vec{\theta} | \mathcal{D}) \approx q(\vec{\theta}) \quad \text{how?}$$

$$p(y | \mathcal{D}) \approx \int p(y | \vec{x}, \vec{\theta}) q(\vec{\theta}) d\vec{\theta}$$

$$q_1(\vec{\theta})$$





$$q_{\lambda^*}(\vec{\theta}) \approx p(\vec{\theta}|\mathcal{D})$$

### 3.6. Laplace approximation

for approximate Bayesian inference

The problem: we know

$f(z)$  - unnormalized density

but we want to know:

$$p(z) = \frac{1}{C} f(z), \quad z \in \mathbb{R} \text{ (or } \mathbb{R}^d)$$

$$1 = \int p(z) dz = \int \frac{1}{C} f(z) dz$$

$$C = \int f(z) dz - \text{but we can't}$$

evaluate the integral analytically

Ex. Bayesian inference

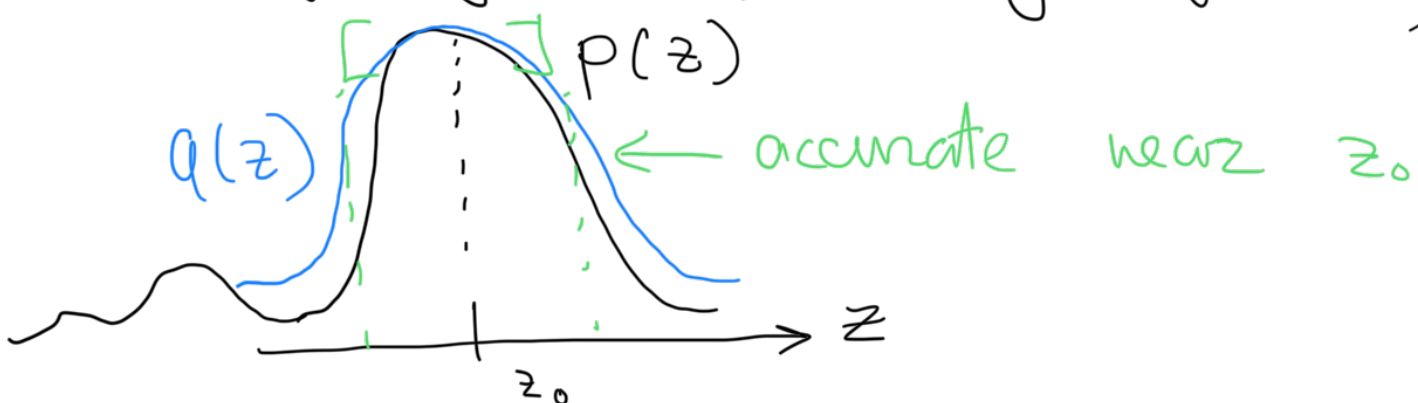
$$p(\theta|\mathcal{D}) = \frac{1}{p(\mathcal{D})} p(\mathcal{D}|\theta) p(\theta)$$

$$f(\theta) = p(\mathcal{D}|\theta) p(\theta), \quad C = p(\mathcal{D})$$

Solution: find  $q(z)$  s.t.

$$q(z) \approx p(z)$$

Idea:  $q(z) \approx p(z)$  for the areas of high (max) density of  $f(z)$



(\*) 1.  $\frac{\partial f(z)}{\partial z} \Big|_{z=z_0} = 0 \leadsto z_0$  - the mode of  $p(z)$

$$2. \quad \ln f(z) \approx \ln f(z_0) - \frac{1}{2} A (z - z_0)^2,$$

- Taylor decomposition for  $\log f(z)$   
around  $z_0$

- No first-order term because of (\*)

$$\bullet \quad A = - \frac{\partial^2}{\partial z^2} \ln f(z) \Big|_{z=z_0}$$

$$f(z) \approx \underbrace{f(z_0)} \exp\left(-\frac{1}{2} \underbrace{A(z-z_0)^2}\right) = q(z)$$

$$C = \sqrt{2\pi} |A|^{-1/2} / f(z_0)$$

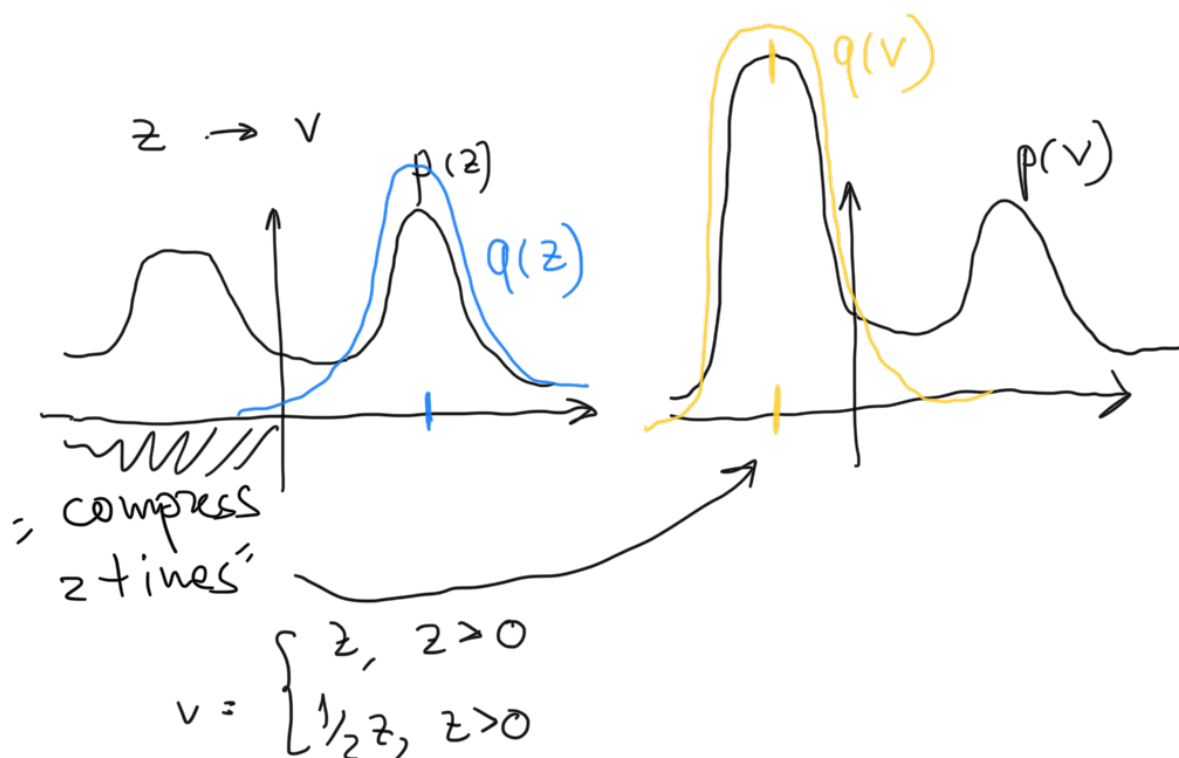
$$q(z) = \mathcal{N}(z | z_0, A^{-1})$$

Why:

BVM theorem states that for larger samples the posterior is close to Gauss.

Why not:

1. local approach
2. Unimodal
3. Depends on parametriz.



### 3.7. Bayesian information criterion (BIC)

$$\begin{aligned}
 Z &= \int_{-\infty}^{+\infty} f(z) dz \approx \int_{-\infty}^{+\infty} f(z_0) \exp\left(-\frac{1}{2} A(z-z_0)^2\right) dz = \\
 &= f(z_0) \int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2} A(z-z_0)^2\right) dz = \\
 &= f(z_0) \frac{(2\pi)^{p/2}}{|A|^{1/2}}, \quad \vec{z} \in \mathbb{R}^p
 \end{aligned}$$

$$\begin{aligned}
 p(\mathcal{D}) &= \int p(\mathcal{D}|\vec{\theta}) p(\vec{\theta}) d\vec{\theta} \\
 f(\vec{\theta}) &= p(\mathcal{D}|\vec{\theta}) p(\vec{\theta})
 \end{aligned}$$

$$Z = p(\mathcal{D})$$

$$\begin{aligned}
 \ln p(\mathcal{D}) &\approx \ln p(\mathcal{D}|\vec{\theta}_{MAP}) + \ln p(\vec{\theta}_{MAP}) + \\
 &\quad + \frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(A)) \approx
 \end{aligned}$$

$$|\mathcal{D}| = n \rightarrow \infty$$

$$\approx \ln p(\mathcal{D}|\vec{\theta}_{MAP}) + 0 + 0 - \frac{1}{2} p \ln(n)$$

A - diagonal D - i.i.d. observations

$$A = \begin{pmatrix} \frac{\partial^2 \ln p(\mathcal{D}|\vec{\theta})}{\partial \theta_1^2} & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \ddots \\ 0 & \dots & \frac{\partial^2 \ln p(\mathcal{D}|\vec{\theta})}{\partial \theta_p^2} & 0 \end{pmatrix}$$

$$\begin{aligned}
 \frac{\partial^2 \ln p(\mathcal{D}|\vec{\theta})}{\partial \theta_i^2} &= \frac{\partial^2 \sum_{j=1}^n \ln p(y_j|\vec{\theta})}{\partial \theta_i^2} \approx n \mathbb{E} \frac{\partial^2 \ln p(y_i|\vec{\theta})}{\partial \theta_i^2} \\
 &\approx n \cdot c
 \end{aligned}$$

$$\ln \det(A) = \ln[(n \cdot c)^p] = p \ln n + c'$$

$$\text{BIC} = p \ln n$$

### 3.8. Bayesian Logistic Regression

$$p(\vec{\theta}) = \mathcal{N}(\vec{\theta} | \vec{\mu}_0, S_0)$$

$$p(\vec{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\vec{\theta}) p(\vec{\theta})$$

$$\begin{aligned}
 \ln p(\vec{\theta}|\mathcal{D}) &\propto \sum_{i=1}^n (y_i \ln p_i + (1-y_i) \ln(1-p_i)) - \\
 &\quad - \frac{1}{2} (\vec{\theta} - \vec{\mu}_0)^T S_0^{-1} (\vec{\theta} - \vec{\mu}_0), \quad p_i = \sigma(\vec{\theta}^T \vec{x}_i)
 \end{aligned}$$



$$p(\vec{\theta} | \mathcal{D}) \approx \underbrace{q(\vec{\theta})}_{S_n^{-1}} = \mathcal{N}(\vec{\theta} | \underbrace{\vec{\theta}_{MAP}}_{\vec{\theta}_{MAP}}, \underbrace{S_n}_{S_n})$$

$$S_n^{-1} = -\nabla \nabla \ln p(\vec{\theta} | \mathcal{D}) = \sum_{i=1}^n p_i(1-p_i) \underbrace{\vec{x}_i \vec{x}_i^T}_{p \times p} + S_0^{-1}$$

## Predictive distribution

$$p(y=0 | \vec{x}, \mathcal{D}) = \int p(y=0 | \vec{x}, \vec{\theta}) p(\vec{\theta} | \mathcal{D}) d\vec{\theta} \approx$$

$$\approx \int \underbrace{p(y=0 | \vec{x}, \vec{\theta})}_{\delta(\vec{x}^T \vec{\theta})} \underbrace{q(\vec{\theta})}_{\text{Gaussian}} d\vec{\theta} \equiv$$

- not analytical ☹

$$\delta(\vec{x}^T \vec{\theta}) = \int \underbrace{\delta(a - \vec{\theta}^T \vec{x})}_{\text{delta function}} \delta(a) da$$

$$\equiv \int \delta(a) p(a) da, \equiv$$

$$\underline{p(a)} = \int \delta(a - \vec{\theta}^T \vec{x}) \underline{q(\vec{\theta})} d\vec{\theta}$$

$$p(a) = \mathcal{N}(a | \mu_a, \sigma_a^2)$$

$$\mu_a = \mathbb{E}a = \int p(a) da = \int q(\vec{\theta}) \vec{\theta}^T \vec{x} d\vec{\theta} = \vec{\theta}_{MAP}^T \vec{x}$$

$$\sigma_a^2 = \text{Var}(a) = \int p(a) \{a^2 - (\mathbb{E}a)^2\} da =$$

$$= \int q(\vec{\theta}) \{(\vec{\theta}^T \vec{x})^2 - (\vec{\theta}_{MAP}^T \vec{x})^2\} d\vec{\theta} =$$

$$= \vec{x}^T \underbrace{S_n^{-1}} \vec{x}$$

$$\equiv \int \delta(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da \approx$$

$$\approx \int \underbrace{\Phi(\lambda a)}_{\text{CDF for standard Gaussian distribution}} \mathcal{N}(a | \mu_a, \sigma_a^2) da =$$

CDF for standard Gaussian distribution

$$= \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right) \approx \boxed{\delta(k(\sigma^2) \mu_a)},$$

$$k(\sigma_a^2) = \sqrt{\frac{1}{1 + \frac{\sigma_a^2}{\lambda^2}}}$$