

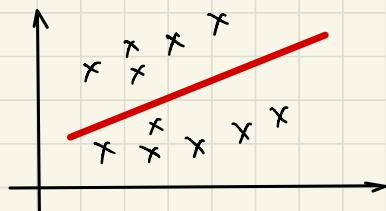
Generalized linear models and exponential families.

Motivation:

- least squares:

$$\min_w \frac{1}{2} \sum_{i=1}^N (y_i - x_i^T w)^2$$

- . prob. model

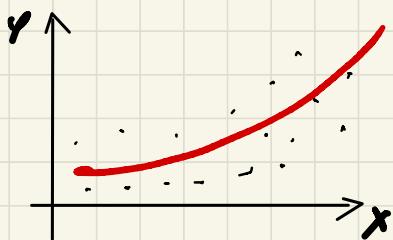


linear model
gaussian noise

$$y = \hat{x}^T w + \hat{\epsilon}; \quad \hat{\epsilon} \sim N(0, \sigma^2)$$
$$y \sim N(x^T w, \sigma^2 I)$$
$$\rightarrow \text{MLE} \rightarrow \hat{w}$$

→ very simple model

In general we will have more complicated data



We can do:

- feature augmentation
 (x, x^2, x^3, \dots)

- Richer class of probabilistic model
(beyond Gaussian and logistic)

→ General linear models
Expon. families.

• Exponential family:

• Logistic regression:

$$\left\{ \begin{array}{l} P(y=1|x,w) = \Gamma(x^T w) = \frac{e^\gamma}{1+e^\gamma} \\ P(y=0|x,w) = 1 - \Gamma(x^T w) = \frac{1}{1+e^\gamma} \end{array} \right. , \quad \gamma = w^T x = x^T w$$

single func.: $P(y|\gamma) = \frac{e^{\gamma y}}{1+e^\gamma}$

$y=1$	$\frac{e^\gamma}{1+e^\gamma}$
$y=0$	$\frac{1}{1+e^\gamma}$

$$P(y|y) = e^{[\gamma y - \log(1+e^\gamma)]}$$

$P(y|\gamma) = e^{\gamma \Phi(y) - A(\gamma)}$

exp. form.

• y is related to $N = \mathbb{E}[y]$ through a non-linear relation:

$\eta = \ln \frac{N}{1-N} \rightarrow N = \Gamma(\eta)$

→ link function

Exponential family:

$$P(y|\eta) = h(y) \exp[\eta^T \phi(y) - A(\eta)]$$

- $y \in \mathbb{R}, \mathbb{R}_+, [0, 1], \mathbb{N}$, can be cont. or disc.
- η : natural parameter
- $h(y) > 0$

$\phi(y)$: sufficient statistic

- $A(\eta)$: cumulant, log partition

$$\int h(y) \exp(\eta^T \phi(y) - A(\eta)) dy = 1$$

$$\Rightarrow A(\eta) = \log \left[\int h(y) \exp(\eta^T \phi(y)) dy \right]$$

$M := \{ \eta ; \int h(y) \exp(\eta^T \phi(y)) dy < \infty \}$

\hookrightarrow natural parameter space

$$\Rightarrow \text{degrees of freedom } (h, \phi, \eta)$$

Ex: $y = h_0, 1$

• Bernoulli dist. with par. ν

$$P(y|\nu) = \nu^y (1-\nu)^{1-y} = \left(\frac{\nu}{1-\nu} \right)^y \cdot (1-\nu)$$

$$= e^{\underbrace{\log\left(\frac{\nu}{1-\nu}\right)y}_{\text{dependency}} + \underbrace{\log(1-\nu)}_{\text{indep. of } y}}$$

$$[\eta^T \phi(y) - A(\eta)] \quad (\nu = \eta^T \eta)$$

$$\phi(y) = y, \quad \eta = \ln \frac{\nu}{1-\nu}; \quad h(y) = 1, \quad A(\eta) = -\log(1-\nu) = \log(1+e^\eta)$$

N is the expected value of y :

$$N = \mathbb{E}(\phi(y))$$

(2 ways to parametrize dist.: N or η)

• Gaussian: $N(\mu, \sigma^2)$

$$P(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

$y \in \mathbb{R}$

• we have two parameters (μ, σ^2)

$$\begin{aligned} P(y | \mu, \sigma^2) &= e^{-\frac{1}{2} \log(2\pi\sigma^2) - \frac{y^2}{2\sigma^2} + \frac{2\mu y}{\sigma^2} - \frac{\mu^2}{2\sigma^2}} \\ &= e^{(\frac{\mu}{\sigma^2}; -\frac{1}{2\sigma^2})(\frac{y}{\sigma^2}) - \frac{\mu^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2)} \\ &= e^{\eta^T \phi(y) - A(\eta)} \end{aligned}$$

$$\eta = \begin{pmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{pmatrix}; \quad h(y) = 1 \quad A(\eta) = \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)$$

$$\phi(y) = \begin{pmatrix} y \\ y^2 \end{pmatrix} \quad = -\frac{\eta_1^2}{4\eta_2} + \frac{1}{2} \log\left(-\frac{\pi}{\eta_2}\right)$$

$$A(\eta) = -\frac{\eta_1^2}{4\eta_2} - \log\left(\frac{\eta_2}{\pi}\right)$$

Link fact: $\eta_1 = \mu/\sigma^2$ $\eta_2 = -1/2\sigma^2$

$$N = -\frac{\eta_1}{2\eta_2}; \quad \sigma^2 = -\frac{1}{2\eta_2}$$

• Poisson with mean λ . $y \in \mathbb{N}$

$$\begin{aligned} P(y|\lambda) &= \frac{\lambda^y}{y!} \cdot e^{-\lambda} \\ &= \frac{1}{y!} \cdot e^{y \ln(\lambda) - \lambda} \\ &= h(y) \cdot e^{\eta^T \phi(y) - A(\eta)} \end{aligned}$$

$$\begin{aligned} h(y) &= \frac{1}{y!} \\ \phi(y) &= y \end{aligned} ; \quad \eta = g(\lambda) = \ln(\lambda) \Rightarrow \lambda = e^\eta$$

⚠ Cauchy dist. is not a member of the exp. family.

Basic properties of $A(\eta)$

- $A(\eta)$ is convex
- $\nabla_\eta A(\eta) = \mathbb{E}[\phi(y)]$
- $\nabla_\eta^2 A(\eta) = \mathbb{E}[\phi(y)\phi(y)^T] - \mathbb{E}[\phi(y)]\mathbb{E}[\phi(y)]^T$

Ex: The case of logistic regression

$$A(\eta) = \log(1 + e^\eta)$$

$$\nabla A(\eta) = \frac{e^\eta}{1 + e^\eta} = \tau(\eta)$$

$$\nabla^2 A(\eta) = \tau(\eta)(1 - \tau(\eta)) \geq 0$$

$$\nabla A(\eta) = \tau(\eta) = \mathbb{E}[\phi(y)] = \mathbb{E}[y]$$

$$\nabla^2 A(\eta) = \tau(\eta)(1 - \tau(\eta)) = N(1 - N) = \text{Var}(y)$$

Link function

so the fact g s.t. $y = g(\mathbb{E}[\phi(y)])$

Parameter estimation

- Fixed family (h, ϕ) + data $(y_i)_{i=1}^N$ i.i.d

→ Recover η

- MLE: $P(y|\eta) = h(y) e^{y^T \phi(y) - A(\eta)}$

$$\begin{aligned} \mathcal{L}(\eta) &= -\ln(P(y|\eta)) \\ &= \sum_{i=1}^N -\ln(P(y_i|\eta)) \\ &= \sum_{i=1}^N [-\ln(h(y)) - y^T \phi(y) + A(\eta)] \end{aligned}$$

$$\nabla \mathcal{L}(\eta) = \sum_{i=1}^N -\phi(y_i) + N \nabla A(\eta)$$

$$= -\sum_{i=1}^N \phi(y_i) + N \mathbb{E}[\phi(y)]$$

$$\nabla \mathcal{L}(\eta) = 0 \Leftrightarrow \mathbb{E}[\phi(y)] = \frac{1}{N} \sum_{i=1}^N \phi(y_i)$$

with a link fct g:

$$y = g(\mathbb{E}[\phi(y)]) \Rightarrow \eta = g\left(\frac{1}{N} \sum_{i=1}^N \phi(y_i)\right)$$

Generalized linear model:

(x, y) iid

$$P(y|w, x) = h(y) \cdot e^{x^T w \phi(y) - A(x^T w)}$$

for $\eta = x^T w$: the linear prediction

$S_{\text{Train}} = (x_i, y_i)_{i=1}^N$, Goal: Find w
How: MLE

Negative log likelihood:

$$\mathcal{L}(w) = - \sum_{i=1}^N \ln(P(y_i|x_i, w))$$

$$= - \left[\sum_{i=1}^N \log(h(y_i)) + x_i^T w \phi(y_i) - A(x_i^T w) \right]$$

\mathcal{L} is convex.

$$\begin{aligned} \nabla \mathcal{L}(w) &= \sum_{i=1}^N -x_i \phi(y_i) + A'(x_i^T w) x_i \\ &= \sum_{i=1}^N -x_i \phi(y_i) + \mathbb{E}[\phi(y_i)] x_i \\ &= \sum_{i=1}^N -x_i \phi(y_i) + g^{-1}(x_i^T w) x_i \end{aligned}$$

$$\nabla \mathcal{L}(w) = 0 \Leftrightarrow g^{-1}(x_i^T w) x_i = \sum_{i=1}^N x_i \phi(y_i)$$

$$\Leftrightarrow g^{-1}(x_i^T w) x_i - \sum_{i=1}^N x_i \phi(y_i) = 0 \quad (\text{LS: } g = \text{id})$$

$$\Leftrightarrow X^T [g^{-1}(Xw) - y] = 0 \quad (\text{LR: } g^{-1} = T)$$

Recap

- Linear Model: $y = x^T w + \epsilon \rightsquigarrow \text{LS. estimator}$
- Logistic regression: $P(y=1|x, w) = \sigma(x^T w)$
- Exponential family: $P(y|\eta) = h(y) e^{\eta^T \phi(y) - A(\eta)}$
 - h, ϕ you can decide
 - η : nature par.
 - A : log-partition (cumulant)
 - $A(\eta)$ is convex
 - $\nabla A(\eta) = E[\phi(y)]$
- G.L.M.

$$P(y|x, w) = h(y) e^{(x^T w \phi(y) - A(x^T w))}$$

\leadsto with NLL find \hat{w}