

Lecture 7 - linear & logistic regression

Thursday, September 26, 2024 2:37 PM

today: linear regression
logistic

(asymptotic) properties MLE:

under suitable regularity conditions on Θ & $p(x; \theta)$

$$\hat{\theta}_n = \underset{\theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta)$$

a) $\hat{\theta}_n \xrightarrow{P} \theta$ 'consistent'

b) CLT
(central limit theorem)

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\text{dist}} N(0, I(\theta)^{1/2})$$

$$[D_n \sim p_{\theta}^{\otimes n}]$$

information matrix

c) asymptotically optimal

i.e. it has minimal asymptotic scaled variance among all "reasonable" estimators

(Cramer-Rao lower bound)

(always true)
↳ d)
(not just asymptotic)

invariance: MLE is preserved under reparametrization

suppose have a bijection $f: \Theta \rightarrow \Theta'$

$$\text{then } \hat{f}(\hat{\theta}) = f(\hat{\theta})$$

$$\text{example: } \hat{(\sigma^2)} = (\hat{\sigma})^2$$

$$\hat{\sin \sigma^2} = \sin \hat{\sigma^2}$$

* if not a bijection, can generalize the MLE with profile likelihood

suppose $g: \Theta \rightarrow \mathcal{L}$

$$\text{profile likelihood} \triangleq L(n) = \max_{\theta: g(n)=g(\theta)} p(\text{data}, \theta)$$

$$\text{define } \hat{\theta}_{MLE} \triangleq \underset{\theta \in \Theta}{\operatorname{argmax}} L(n)$$

then we have

$$\hat{\theta}_{MLE} = g(\hat{\theta}_{MLE})$$

"plug in estimate"

$$N(\mu, \sigma^2)$$

$$\text{e.g. } g(\mu) = \mu^2$$

prediction

want to learn a prediction fct, $h: X \rightarrow \mathcal{Y}$

$$x \in \mathbb{R}^d$$

$$\mathcal{Y} = \{0, 1\} \rightarrow \text{binary classification}$$

$$\mathcal{Y} = \{0, 1, \dots, K-1\} \rightarrow \text{multiclass classification}$$



'prediction model' $\mathcal{Y} = \mathbb{R} \rightarrow$ regression

$$p(x, y) = \underbrace{p(y|x)}_{\text{marginal model on } X} p(x)$$

$p(x, y) = \underbrace{p(y|x)}_{\text{"prediction model"} - \text{marginal model on } X} \underbrace{p(x)}_{\text{regression}}$
 $= \underbrace{p(x|y)}_{\text{"class conditional"}} \underbrace{p(y)}_{\text{prior over class}}$

"generative perspective" (in context of classification) \rightarrow model $p(x)$ as well
vs.

"conditional perspective" \rightarrow only model $p(y|x)$
"more discriminative" \uparrow traditionally called "discriminative"

generative model $p_\theta(x, y)$ MLE	conditional model $p_\theta(y x)$ max conditional likelihood (MCL)	"fully discriminative" model $h_\theta: X \rightarrow Y$ (not nec. derived from $p(y x)$) reg. ERM; etc. $\frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}; h_\theta(x^{(i)})) + \text{reg}(\theta)$
more assumptions \Rightarrow less robust predictions		less assumption more robust

$$\hat{h}(x) \triangleq \underset{\tilde{y} \in Y}{\text{argmin}} \sum_y \hat{p}(y|x) \ell(y, \tilde{y})$$

if $\ell(y, \tilde{y}) = \mathbb{1}\{y \neq \tilde{y}\}$ (0-1 loss) then $\hat{h}(x) = \underset{\tilde{y} \in Y}{\text{argmax}} p(\tilde{y}|x)$

linear regression: derive/motivate with conditional approach to regression ($y \in \mathbb{R}$)

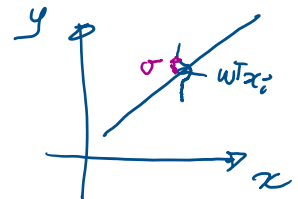
$$p(y|x; w) = N(y | \underbrace{\langle w, x \rangle}_{w^T x}, \sigma^2)$$

$x \in \mathbb{R}^d$
 $w \in \mathbb{R}^d$

parameter

equivalently: $y_i = w^T x_i + \epsilon_i$

where $\epsilon_i | x_i \stackrel{iid}{\sim} N(0, \sigma^2)$



[aside: we use "offset" notation for x

i.e. $x = \begin{pmatrix} \tilde{x} \\ 1 \end{pmatrix}$ $\tilde{x} \in \mathbb{R}^{d-1}$
"constant feature"

thus $\langle w, x \rangle = \langle w_{1:d}, \tilde{x} \rangle + w_d$
["bias" / "offset" (usually denoted as b)]

* dataset $(x_i, y_i)_{i=1}^n$
 $x_i \sim \text{whatever (don't care)}$
 $y_i | x_i \stackrel{iid}{\sim} N(w^T x_i, \sigma^2)$

conditional likelihood:

$$y_i | x_i \sim N(w^T x_i, \sigma^2)$$

conditional likelihood:

$$p(y_{1:n} | x_{1:n}) = \prod_{i=1}^n p(y_i | x_i)$$

$$\log \left[\right] = \sum_{i=1}^n \left[-\frac{(y_i - w^T x_i)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]$$

this is not a concave fct. of σ^2
 $f(u) = -\frac{a}{u} - b \log u$

$$\sigma^2 = \frac{1}{\lambda} \quad \lambda = \frac{1}{\sigma^2}$$

$$f(\lambda) = -a\lambda + b \log \lambda$$

concave concave

$$\frac{\partial}{\partial \lambda} \left(\right) = 0 \Rightarrow \sum_{i=1}^n \left(-\frac{(y_i - w^T x_i)^2}{2} + \frac{1}{2} \frac{1}{\lambda} \right) = 0$$

$$\Rightarrow \hat{\lambda}_{MLE} = \left[\frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2 \right]^{-1}$$

(global max by concavity)

(by invariance of MLE) $\Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$

\Rightarrow conclude that this correct global in σ^2 for w fixed

15h39

"design matrix"

(matrix??)

$$X \triangleq \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix} \quad y \triangleq \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$n \times d$ matrix

$$Xw = \begin{pmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{pmatrix} \in \mathbb{R}^{n \times 1}$$

$$\sum_{i=1}^n (y_i - w^T x_i)^2 = \|y - Xw\|_2^2$$

can rewrite $-\log p(y_{1:n} | X) = \frac{\|y - Xw\|_2^2}{2\sigma^2} + \text{fct.}(\sigma^2)$

design matrix

MCL $\rightarrow \min_w \|y - Xw\|_2^2 \Leftrightarrow$ projection of y on the column space of design matrix X



$$\hat{w}_{MLE} = \arg \min_{w \in \mathbb{R}^d} \|y - Xw\|_2^2$$

"least square"

$$Xw = \sum_{j=1}^d X_{:,j} w_j$$

j^{th} column of X

algebra: want $\nabla_w \rightarrow 0$

$$\frac{\partial}{\partial w} \left[(y - Xw)^T (y - Xw) \right] \stackrel{\text{want}}{=} 0$$

$$\frac{\partial}{\partial w} \left[\|y\|^2 - 2y^T Xw + w^T X^T Xw \right]$$

vector

$$\nabla_w (w^T A w) = (A + A^T) w$$

$$\frac{\partial}{\partial w} [\|y\|^2 - 2y^T Xw + w^T X^T X w]$$

$$0 - 2X^T y + 2X^T X w = 0$$

$$\Rightarrow \boxed{(X^T X)w = X^T y}$$

$$= (A + A^T)w$$

$\| \cdot \|^2$ convex fct. of w
 \Rightarrow stat. pt. are global min

"normal equation"

a) if $X^T X$ is invertible, then have unique sol'n

$$\hat{w}_{MLE} = \underbrace{(X^T X)^{-1}}_{X^+} X^T y$$

prediction on training set:

$$\hat{y} = X \hat{w} = X \underbrace{(X^T X)^{-1} X^T}_{\text{projection on column space of } X} y$$

projection on column space of X

$$X \text{ is } n \times d \Rightarrow \text{rank}(X) \leq \min\{n, d\}$$

$$\text{rank}(X^T X) = \text{rank}(X)$$

$$X^T X \text{ is invertible} \Rightarrow \min\{n, d\} = n \geq d$$

\rightarrow recall geometric perspective?

⊛ if $n < d$ (i.e. high dimension or low data regime)
 then $X^T X$ is not invertible

b) if $(X^T X)$ is not invertible \rightarrow there is no unique sol'n

any \hat{w} s.t. $(X^T X)\hat{w} = X^T y$ is a MLE sol'n

could choose $\hat{w} = \underset{w \text{ s.t. } (X^T X)w = X^T y}{\arg \min \|w\|_2} = X^+ y$ — Moore-Penrose pseudo-inverse

$$X^+ = (X^T X)^{-1} X^T \text{ when } X \text{ is full rank}$$

SVD

$$X = U \Sigma V^T$$

$\downarrow \quad \quad \downarrow \quad \quad \downarrow$
 $n \times d \quad \quad n \times n \quad \quad n \times d$

$$U^T U = I_n \quad V^T V = I_d$$

$$\Sigma = \begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_d & \\ & & & 0 \end{pmatrix}$$

$$X^T = V \Sigma^T U^T$$

$$\Sigma^+ = \begin{pmatrix} \sigma_1^+ & & 0 \\ & \ddots & \\ 0 & & \sigma_d^+ & \\ & & & 0 \end{pmatrix}$$

$$\sigma_i^+ = \begin{cases} 1/\sigma_i & \text{if } \sigma_i \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect

$$\hat{w}_{MAP}(\lambda) \xrightarrow{\lambda \rightarrow 0} \hat{w}_{\text{pseudo-inverse}}$$

regularization: (can be motivated from MAP pt. of new)

suppose we put a prior $p(w) = N(w|0, \sigma^2 I)$

$\sigma^2 I$ is identity matrix

σ^2 "precision coefficient"

log posterior: $\log p(w|\text{data}) = \log p(y_{1:n}|X, w) + \log p(w) + \text{const.}$

precision: $\frac{1}{\sigma^2}$

$$\begin{aligned} \text{log posterior: } \log p(w|\text{data}) &= \log p(y_{1:n}|X, w) + \log p(w) + \text{cst.} \\ &= \underbrace{-\frac{1}{2\sigma^2} \|y - Xw\|^2}_{\text{likelihood}} - \underbrace{\frac{\lambda}{2\sigma^2} \|w\|^2}_{\text{prior}} + \text{cst.} \end{aligned}$$

(of w)

MAP here

$$\hat{w}_{\text{MAP}} = \underset{w \in \mathbb{R}^d}{\text{argmin}} \left(\frac{1}{2} \|y - Xw\|^2 + \frac{\lambda}{2} \|w\|^2 \right)$$

"ridge regression"

→ same as "regularized ERM"

with squared loss $\ell(y_i, w^T x_i) = \frac{1}{2} (y_i - w^T x_i)^2$

$$\frac{1}{n} \sum \underbrace{\frac{1}{2} (y_i - w^T x_i)^2}_{\text{squared loss}} + \underbrace{\frac{\lambda}{2n} \|w\|^2}_{\text{regularization}}$$

this obj. is strongly convex in w
 \Rightarrow a unique soln

$$\left[f(\cdot) \text{ is strongly convex} \Leftrightarrow f(\cdot) - \frac{\lambda}{2} \|\cdot\|^2 \text{ is convex in } (\cdot) \right]$$

$$\nabla_w = 0 \Rightarrow \underbrace{(X^T X + \lambda I)}_{\text{always invertible } \lambda > 0} w = X^T y$$

$$\hat{w}_{\text{MAP}} = (X^T X + \lambda I)^{-1} X^T y$$

no problem for $d > n$

• note about σ^2 being a global max

(aside: showing that the σ^2 above is the **global max** is subtle because the objective is not concave in σ^2 . I give more info here for your curiosity, but it is not required for the assignment. And it's easier instead to just use $\lambda = 1/\sigma^2$ instead as the parameterization as I have done, as the objective is then concave in λ ...)

- Formally, to find a global max of a *differentiable objective*, you need to check all **stationary points** (zero gradient points), **as well as the values at the boundary of the domain**.

Thus here, you would need to show that the objective cannot take higher value anywhere at the boundary of the domain (which is the case here (exercise!), as the objective goes to $-\infty$ at the boundary), so you are done (this is the only possible global optimum -- a maximum here, as it should be, given that there are no other stationary points and all values are lower at the boundary, but one could also explicitly check the Hessian to see that it is strictly negative definite at the stationary point, i.e. it looks like a local maximum).

Note that we will see later in the class that the Gaussian is in the exponential family, with a log-concave likelihood in the right ("natural") parameterization, and thus using the invariance principle of the MLE, we could also easily deduce the MLE in the "moment" parameterization which is the usual (μ, σ^2) one, without having to worry about local optima...

- for a cute counter-example illustrating that a differentiable function could have only one stationary point which is a local min but *not a global min* (and thus why one needs to look at the values at the boundary), see:
 - https://en.wikipedia.org/wiki/Maxima_and_minima#Functions_of_more_than_one_variable
 - i.e.

$$f(x, y) = x^2 + y^2(1 - x)^3, \quad x, y \in \mathbb{R},$$

shows. Its only critical point is at $(0,0)$, which is a local minimum with $f(0,0) = 0$. However, it cannot be a global one, because $f(2,3) = -5$.

(see picture of function [here](#))

(and note that the "[Mountain pass theorem](#)" which basically says that if you have a strict local optimum with another point somewhere with the same value, then there must be a saddle point somewhere (a "mountain pass") i.e. another stationary point, **does not hold for this counter-example** as one of the required regularity condition, the "Palais-Smale compactness condition" fails. Here, the saddle point (which should intuitively exist) "happens at infinity", which is why it only has one stationary point despite $(0,0)$ not being a global minimum)

- the moral of the story: intuitions for multivariate optimization are often misleading! (this counter-example would not work in 1d because of [Rolle's theorem](#))