

today: linear regression
logistic 11

[see lecture last year:]

(asymptotic)
properties of MLE:

under suitable regularity conditions on $\Theta \ni p(x; \theta)$ $\hat{\theta}_n = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta)$

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

b) CLT (central limit theorem) $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I(\theta)^{-1})$ $[D_n \sim P_{\theta}^{\otimes n}]$
information matrix

c) asymptotically optimal

(Cramer-Rao lower bound)

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

d) invariance: MLE is preserved under reparameterization

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

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

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

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

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

suppose $g: \Theta \rightarrow \Lambda$

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

define $\hat{\eta}_{MLE} \triangleq \operatorname{argmax}_{\eta \in \Lambda} L(\eta)$

then we have $\boxed{\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})}$

"plug in estimator"

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

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

prediction

want to learn a prediction fct. $h: X \rightarrow Y$

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

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

$Y = \{0, \dots, k-1\} \rightarrow$ multiclass "

$Y = \mathbb{R} \rightarrow$ regression



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

$$= \underbrace{p(x|y)}_{\text{"class conditional"}} \underbrace{p(y)}_{\text{prior over class}}$$

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

vs.

"conditional perspective" \rightarrow only model $p(y|x)$

"more discriminative"

traditionally called "discriminative"

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

$$\hat{h}(x) \triangleq \arg \min_{\tilde{y} \in Y} \sum_y p_{\theta}(y|x) \ell(y, \tilde{y})$$

if $\ell(y, \tilde{y}) = 1\{y \neq \tilde{y}\}$ then $\hat{h}(x) = \arg \max_{\tilde{y} \in Y} p_{\theta}(\tilde{y}|x)$
(0-1 loss)

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

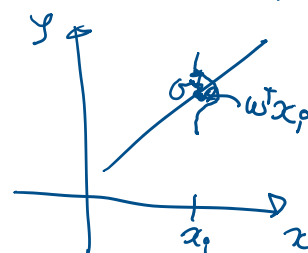
$$p(y|x; w) = N(y | \underbrace{\langle w, x \rangle}_{w^T x}, \sigma^2)$$

parameter

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

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

equivalently: $Y_i = w^T X_i + \epsilon_i$



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

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

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

$$\begin{array}{c} \text{---} \\ | \\ x_i \end{array} \rightarrow x$$

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

[aside: we use "offset" notation for x]

]

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

"constant feature"

$$\text{thus } \langle w, x \rangle = \langle w_{1:d-1}, \tilde{x} \rangle + w_d$$

"bias" / "offset" (usually denoted as b)

* dataset $(x_i, y_i)_{i=1}^n$

$$X_i \sim \text{whatever (don't care)}$$

$$y_i | x_i \sim \text{iid } N(w^T x_i, \sigma^2)$$

conditional likelihood:

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

this is not a concave fct. of σ^2

$$g(u) = \frac{a}{u} + b \log u$$

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

$$\frac{\partial}{\partial \sigma^2} \left(\sum_{i=1}^n \left[-\frac{(y_i - w^T x_i)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] \right) \stackrel{!}{=} 0 \Rightarrow \sum_{i=1}^n \left[\frac{(y_i - w^T x_i)^2}{2(\sigma^2)^2} - \frac{1}{2} \frac{1}{\sigma^2} \right] = 0$$

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

$$\log \rightarrow -\infty \text{ as } \sigma \rightarrow 0 \quad \sigma \rightarrow +\infty$$

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

(see also notes [below](#))

"design matrix" X $\begin{pmatrix} \text{matrix} \end{pmatrix}$ $\begin{pmatrix} \text{vector} \end{pmatrix}$

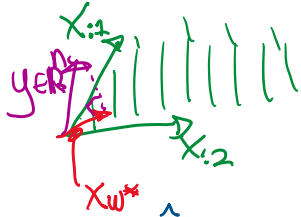
$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$Xw = \begin{pmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{pmatrix} \in \mathbb{R}^{n \times 1} \quad \sum_{i=1}^n (y_i - w^T x_i)^2 = \|y - Xw\|_2^2$$

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

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

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



$$\hat{w}_{MLE} = \underset{w \in \mathbb{R}^d}{\text{argmin}} \|y - Xw\|_2^2 \quad \text{"least square"}$$

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

\uparrow j^{th} column of X

15h32

algebra: want $\nabla_w \xrightarrow{\text{set to}} 0$

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

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

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

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

vector

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

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

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

"normal equation"

a) if $X^T X$ is invertible, then have unique solution
 $d \times d$

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

prediction on training set \hat{w}_{MLE}

$$\hat{y} = X \hat{w} = X (X^T X)^{-1} X^T y$$

projection on column space of X

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

$X^T X$ is invertible $\Leftrightarrow \min\{n, d\} \Rightarrow \boxed{n \geq d}$

(recall geometric perspective)

(*) if $n < d$ (ie. high dimension or low data regime)

then $X^T X$ is not invertible

[MLE $\xrightarrow{\text{replace with above}} \text{MCL}$]

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

any \hat{w} st. $(X^T X) \hat{w} = X^T y$ is a MCL solution

could choose $\hat{w} = \underset{w}{\text{argmin}} \|w\|_2$

$$\text{st. } (X^T X) w = X^T y$$

$= X^+ y$ Moore-Penrose pseudo-inverse

lower bound w - argmin $\|w\|_2$

sol. $(X^T X)w = X^T y$

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

SVD

$X = U \Sigma V^T$
 $n \times d$ $d \times d$ $d \times d$

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

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

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

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

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

problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect.

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

regularization: (can be motivated from MAP point of view)
 suppose we put a prior $p(w) = N(w | 0, \frac{1}{\lambda} I)$ $d \times d$ identity matrix
 λ "precision parameter"

log posterior: $\log p(w | \text{data}) = \log p(y_{1:n} | X, w) + \log p(w) + \text{cst.}$
 $= -\frac{1}{2\sigma^2} \|y - Xw\|^2 - \frac{\lambda}{2} \|w\|^2 + \text{cst.}$ (now)

MAP here $\hat{w}_{MAP} = \arg \min_{w \in \mathbb{R}^d} \frac{1}{2} \|y - Xw\|^2 + \frac{\lambda}{2} \|w\|_2^2$ "ridge regression"

→ same as "regularized ERM" $\rightarrow \frac{1}{n} \sum_{i=1}^n \ell(y_i; h_w(x_i)) + \frac{\lambda}{2} \|w\|^2$
 with squared loss $\ell(y_i; \tilde{w} x_i) = \frac{1}{2} (y_i - \tilde{w} x_i)^2$ Squared loss $\frac{\partial}{\partial n}$ regularization

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

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

$\nabla_w = 0 \Rightarrow (X^T X + \lambda I) w = X^T y$

always invertible $\lambda > 0$

$\hat{w}_{MAP} = (X^T X + \lambda I)^{-1} X^T y$ no problem for

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

or
ridge regression

no problem for $d > n$

good practice: to either standardize features (i.e. make each feature zero mean & unit empirical variance) or normalize \rightarrow make x_i unit norm $\|x_i\|_2 = 1$ or rescale features to $[0,1]$ or $[-1,1]$

- **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.)

- 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](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 conditions, 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](#))