

Lecture 7 - scribbles

Tuesday, September 24, 2019 14:26

today : linear regression
logistic regression

Prediction

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

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

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

$\{0, 1, \dots, K\} \rightarrow$ multi class "



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

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

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

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

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

"more discriminative"

\rightarrow traditionally called "discriminative"

generative	conditional	"fully discriminative"
model $p_\theta(x, y)$	model $p_\theta(y x)$	model $h_\theta: X \rightarrow Y$
MLE	max. conditional likelihood	(not nec. derived from $p(y x)$)
more assumptions \Rightarrow less robust for prediction		less assumptions more robust

Linear regression : derive / motivate with conditional approach to regression (Yeik)

$$p(y|x; \omega) = N(y | \underbrace{\omega^T x}_{\text{parameter } \omega \in \mathbb{R}^d, x \in \mathbb{R}^d}, \sigma^2)$$

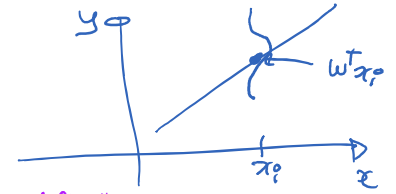
$$N(\mu, \sigma^2)$$

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

equivalently: $y_i = w^T x_i + \epsilon_i$ where $\epsilon_i | x_i \sim \text{iid } N(0, \sigma^2)$

[aside: we'll use "offset" notation for x

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



"bias/offset"

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

• dataset $(x_i, y_i)_{i=1}^n$ $x_i \sim \text{whatever}$

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

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

not a concave fct. of σ^2 opt?

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

$\frac{\partial}{\partial \sigma^2}(\quad) = 0 \rightarrow \sum_{i=1}^n \left[-\frac{(y_i - w^T x_i)^2}{2(\sigma^2)^2} - \frac{1}{2} \frac{1}{\sigma^2} \right] = 0$

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

[global max since obj. $\rightarrow -\infty$ at boundary $\sigma^2 = 0$ or $\sigma^2 \rightarrow \infty$]

15h34

$\left[-\frac{a}{x} - \log x \right]$

$\frac{d}{dx} \begin{bmatrix} -\frac{a}{x} - \log x \end{bmatrix} = \frac{a}{x^2} - \frac{1}{x}$

$\frac{d}{dx^2} = -\frac{2a}{x^3} + \frac{1}{x^2}$ \leftarrow get sign switch \Rightarrow not concave

"design matrix" $X \triangleq$ $n \times d$ matrix

$\begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$

vector $y \triangleq \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

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

can rewrite $-\log p(y:n | X) = \frac{\|y - Xw\|^2}{2\sigma^2} + \text{function}(w)$

↑
design matrix

MCL \rightarrow minimizing $\|y - Xw\|^2 \Leftrightarrow$ projecting y on the column space of design matrix X

(geometric pt. of view)



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

↑
 j^{th} column of X

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

algebra: want $\nabla_w = 0$

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

$$\frac{\partial}{\partial w} (\|y\|^2 - 2y^T Xw + w^T X^T X w) = 0$$

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

$$\Rightarrow \boxed{X^T X w^* = X^T y} \quad \text{"normal equation"}$$

vector

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

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

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

$$X^T X \xrightarrow{\text{need}} d \times d \text{ matrix}$$

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

$$X^T X \text{ invertible} \Rightarrow \boxed{n \geq d}$$

prediction
on training
set

$$\hat{y} = X \hat{w} = X \underbrace{(X^T X)^{-1} X^T}_{\text{projection operator}} y$$

on column space of X (recall geometric)

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

b) what if $X^T X$ is not invertible?

any \hat{w} s.t. $(X^T X) \hat{w} = X^T y$ is a MLE estimator

could choose $\hat{w} = \arg \min_{w: X^T X w = X^T y} \|w\|$ \leftarrow Moore-Penrose pseudo-inverse $X^+ = (X^T X)^+ X^T$ when X is full rank

problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect

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

suppose we put prior $p(w) = N(w|0, \frac{1}{\lambda} I)$ \leftarrow "precision" parameter

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

MAP here

$$\hat{w}_{\text{MAP}} = \arg \min_w \frac{1}{2} \|y - Xw\|^2 + \frac{\lambda}{2} \|w\|^2$$

"ridge regression"

same as "regularized" ERM

$$\frac{1}{n} \sum_{i=1}^n \underbrace{\ell(y_i, h_w(x_i))}_{\text{empirical error}} + \underbrace{\frac{\lambda}{2n} \|w\|^2}_{\text{regularization}}$$

X is "strongly convex" in w

$$\nabla_w = 0$$

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

always invertible for $\lambda > 0$

$f(\cdot)$ is λ -strongly convex

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

\Rightarrow unique solution

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

no problem for $d > n$

$$\hat{w}_{\text{map or ridge regression}} = (X^T X + \lambda I)^{-1} X^T y \quad \text{no problem for } d > n$$

- **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.)

- 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 need 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](#))