

today: linear regression
logistic "

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

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

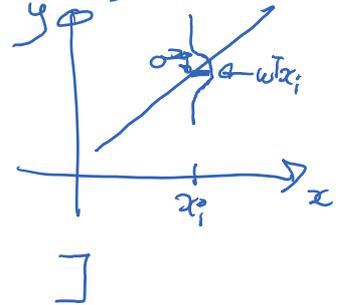
$$\begin{matrix} x \in \mathbb{R}^d \\ w \in \mathbb{R}^d \end{matrix}$$

equivalently: $Y_i = w^T X_i + \epsilon_i$ where $\epsilon_i | X_i \sim N(0, \sigma^2)$

[aside: we use "offset" notation for x

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

$$\text{thus } \langle w, x \rangle = \langle w_{1:d-1}, \tilde{x} \rangle + \underbrace{w_d}_{\text{"bias" / "offset"}} \quad \downarrow$$



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

conditional likelihood

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

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

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

$$\Rightarrow \frac{\partial}{\partial \sigma^2} \log L = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2$$

obj $\rightarrow -\infty$ as $\sigma \rightarrow 0$
or $\sigma \rightarrow +\infty$

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

"Hessian matrix" $\nabla^2 L = \frac{1}{\sigma^4} \sum_{i=1}^n (y_i - w^T x_i)^2 x_i x_i^T$

(see also

"design matrix" $X \triangleq$ $\begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix}$

matrix $\begin{pmatrix} \square & \square \\ \square & \square \end{pmatrix}$
 $n \times d$ matrix

for w fixed

(see also note below)

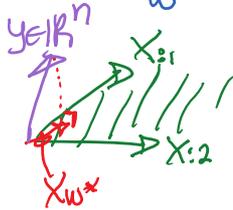
vector $y \triangleq \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

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

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

can rewrite $-\log p(y_{1:n} | X) = \underbrace{\|y - Xw\|_2^2}_{\text{design matrix}} + \text{const}$

MLE $\rightarrow \min_w \|y - Xw\|_2^2 \iff$ projecting y on the column space of design matrix X



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

j^{th} column of X

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

algebra: want ∇_w set to 0

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

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

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

(convex fct. of $w \rightarrow$ stat. pt. is global min)

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

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

"normal equation"

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

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

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

prediction on training set $\hat{y} = X \hat{w} = X (X^T X)^{-1} X^T y$

$X^T X$ is invertible $\Rightarrow n \geq d$

projection matrix on column space of X

(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 estimate

could choose $\hat{w} = \underset{w: (X^T X)w = X^T y}{\operatorname{argmin}} \|w\|_2 = X^+ y$ Moore-Penrose pseudo-inverse (see Wikipedia)

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

$n \times d$ $n \times d$

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

$d \times n$ $d \times n$

$$\Sigma = \begin{pmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ & & \sigma_d & \\ & & & 0 \end{pmatrix} \quad \Sigma^T = \begin{pmatrix} \sigma_1^+ & & & 0 \\ & \ddots & & \\ & & \sigma_d^+ & \\ & & & 0 \end{pmatrix} \quad \sigma_i^+ = \begin{cases} \frac{1}{\sigma_i} & \text{if } \sigma_i > 0 \\ 0 & \text{o.w.} \end{cases}$$

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

problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect

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

16h00

regularization: (can be motivated from MAP point of view)

Suppose we put a prior $p(w) = N(w | 0, \frac{\sigma^2}{\lambda} I)$

λ ← 'precision' parameter
 I ← $d \times d$ identity matrix

$$\begin{aligned} \log \text{posterior: } \log p(w | \text{data}) &= \log p(y_{1:n} | X, w) + \log p(w) + \text{const.} \\ &= \frac{1}{2\sigma^2} \|y - Xw\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 + \text{const.} \end{aligned}$$

MAP here

$$\hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} \|y - Xw\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 \right]$$

"ridge regression"

same as "regularized" ERM

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, w^T x_i) + \frac{\lambda}{2} \|w\|_2^2$$

empirical error
regularization

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

this is obj. is strongly convex in $w \Rightarrow$ a unique sol'n

$f(\cdot)$ is λ strongly convex

$\Leftrightarrow f(\cdot) = \frac{\lambda \|\cdot\|_2^2}{2}$ is convex in (\cdot)

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

always invertible for $\lambda > 0$

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

ridge regression

no problem for $d > n$

good practice: to either standardize features i.e. make each feature zero mean
or
normalize \leftarrow make x_i unit norm $\|x_i\|_2 = 1$ $\left\{ \begin{array}{l} \text{unit empirical variance} \\ \text{or} \\ \text{scale features to } [0,1] \text{ or } [-1,1] \end{array} \right.$

Logistic regression

setup: binary classification $y = \{0,1\}$ $X \in \mathbb{R}^d$

generative model motivation:

suppose only assumption is there exists a pdf (densities) in \mathbb{R}^d for each class conditionals

$$p(x|Y=1) \neq p(x|Y=0)$$

$$p(Y=1|X=x) = \frac{p(Y=1, X=x)}{p(Y=1, X=x) + p(Y=0, X=x)} \} p(X=x)$$

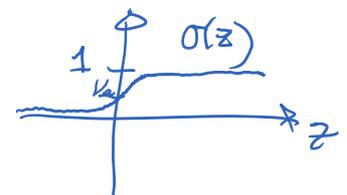
$$\stackrel{!}{=} \frac{1}{1 + \frac{p(Y=0, X=x)}{p(Y=1, X=x)}} = \frac{1}{1 + \exp(-f(x))}$$

where $f(x) \triangleq \log \frac{p(X=x|Y=1)}{p(X=x|Y=0)} + \log \frac{p(Y=1)}{p(Y=0)}$
 "log odds" \leftarrow class-conditional ratio \leftarrow prior odd ratio

in general,

$$p(Y=1|X=x) = \sigma(f(x))$$

where $\sigma(z) \triangleq \frac{1}{1 + \exp(-z)}$
 "sigmoid function"



some properties of $\sigma(z)$:

$$\sigma(-z) = 1 - \sigma(z) \quad [\sigma(z) + \sigma(-z) = 1]$$

$$\frac{d}{dz} \sigma(z) = \sigma(z) \sigma(-z) = \sigma(z) (1 - \sigma(z))$$

\Rightarrow to motivate linear logistic regression, consider class conditionals

in the exponential family

$$p(x|\eta) \triangleq h(x) \exp(\eta^T T(x) - A(\eta))$$

"canonical parameter"
 "sufficient statistics"
 "scalar fct. : log partition fct. normalized"
 these specify the "flat" exponential family that we are considering

Gaussian: $\log p(x|\mu, \sigma^2) = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{(x-\mu)^2}{2\sigma^2}$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = \exp\left[-\frac{x^2}{2\sigma^2} - x\frac{\mu}{\sigma^2} + \frac{\mu^2}{2\sigma^2}\right]$$

let $T(x) = \begin{bmatrix} -x^2/2 \\ x \end{bmatrix}$

$\eta(\mu, \sigma^2) = \begin{bmatrix} 1/2\sigma^2 \\ \mu/\sigma^2 \end{bmatrix}$

$A(\eta) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{\mu^2}{2\sigma^2}$

$$p(x|Y=1) = p(x|\eta_1)$$

$$p(x|Y=0) = p(x|\eta_0)$$

log odds $f(x) = \log \frac{p(x|Y=1)}{p(x|Y=0)} + \log \frac{p(Y=1)}{p(Y=0)}$

$$= (\eta_1 - \eta_0)^T T(x) + A(\eta_0) - A(\eta_1) + \log \frac{\pi}{1-\pi}$$

$$\triangleq w^T \phi(x)$$

where $w = \begin{pmatrix} \eta_1 - \eta_0 \\ A(\eta_0) - A(\eta_1) + \log \frac{\pi}{1-\pi} \end{pmatrix}$ $\phi(x) = \begin{pmatrix} T(x) \\ 1 \end{pmatrix}$

get logistic regression model

$$p_w(Y=1|X=x) = \sigma(w^T \phi(x))$$

"feature map"

decision boundary $\{x \mid w^T \phi(x) > 0\}$

exercise to reader: try argument above with $p(x|y) = N(x|\mu_y, \Sigma_y)$

• if $\Sigma_0 = \Sigma_1$, then $\phi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix}$

• otherwise $\phi(x) = \begin{pmatrix} -xx^T \\ x \\ 1 \end{pmatrix}$

• otherwise $\ell(x) = \left(\begin{array}{c} -xx^T \\ x \\ 1 \end{array} \right)^T \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

- **note about sigma² being a global max**

(**aside:** showing that the sigma² above is the **global max** is subtle because the objective is not concave in sigma². 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 -infinity 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 (mu, sigma²) 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](#))