

today: linear regression
logistic regression

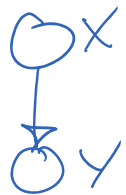
Prediction:

want learn a prediction fct. $h: X \rightarrow \mathcal{Y}$
 $x \in \mathbb{R}^d$

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

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

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



"prediction model" model over X

$$P(x,y) = \overbrace{P(y|x)} \overbrace{P(x)}$$

prior over classes

$$= \underbrace{P(x|y)} \overbrace{P(y)}$$

"class conditional"

"generative perspective" \rightarrow model $p(x)$ as well

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

\uparrow traditionally called "discriminative"

more disc. \rightarrow

... gen | conditional | "fully discriminative"

more use. →

gen	conditional	"fully discriminative"
model $p_{\theta}(x, y)$	model $p_{\theta}(y x)$	model $h_{\theta}: X \rightarrow Y$ (not nec. $p(y x)$)
MLE	max conditional likelihood	surrogate loss minimization ↳ use $l(y, \hat{y})$ for estimation
more assumptions ⇒ less robust for prediction		examples: SVM more robust

Linear regression: conditional approach to regression ($Y \in \mathbb{R}$)

$$p(y|x; w) = N(y | \langle w, x \rangle, \sigma^2)$$

parameter

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

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

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

[aside: we'll use "offset" notation for x i.e. $x = \begin{pmatrix} \tilde{x} \\ 1 \end{pmatrix}$ $\tilde{x} \in \mathbb{R}^{d-1}$ "constant features" "bias/offset"]

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

• dataset $(x_i, y_i)_{i=1}^n$

$X_i \sim \text{whatever}$

$$| \quad \Delta \rightarrow + \|u - x_{..}\|^2 - | |$$

$x_i \sim \text{whatever}$

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

$\frac{\partial}{\partial \sigma^2} \rightarrow + \frac{\|y - Xw\|^2}{2\sigma^4} - \frac{1}{2\sigma^2}$
 $\frac{\partial^2}{\partial (\sigma^2)^2} \rightarrow - \frac{\|y - Xw\|^2}{2\sigma^6} + \frac{1}{2\sigma^4}$

flip sign when σ^2 varies
 ↓
not concave

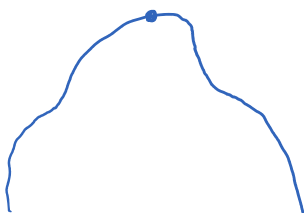
conditional likelihood $p(y_{1:n} | x_{1:n}) = \prod_{i=1}^n p(y_i | x_i)$

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

$\frac{\partial}{\partial \sigma^2} (\dots) = 0$

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

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



(see [note below](#) about this global max)

"design matrix"

matrix X
 $n \times d$ matrix

$$\begin{pmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{pmatrix} \in \mathbb{R}^{n \times d}$$

vector

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

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

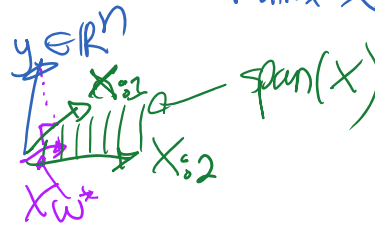
$\|y - Xw\|_2^2 = \sum_{i=1}^n (y_i - w^T x_i)^2$

can rewrite $-\log p(y_{i:n}|X) = \underbrace{\|y-Xw\|^2}_{\text{RSS}} + \text{function}(\sigma^2)$

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

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

\uparrow
j-th col. of X



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

algebra: want $\nabla_w = 0$

$$\frac{\partial}{\partial w} \left((y-Xw)^T (y-Xw) \right) = 0$$

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

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

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

"normal equation"

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

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

$X^T X \rightarrow d \times d$ matrix
 $d \times n \quad n \times d$

$$\hat{\omega}_{MSE} = (X^T X)^{-1} X^T y$$

$X^T X \rightarrow d \times d$ matrix
 $d \times n \quad n \times d$

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

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

prediction on training set $\hat{y} = X \hat{\omega} = \underbrace{X(X^T X)^{-1} X^T}_{\text{projection operator on column space of } X} y$

(see back geometric viewpoint)

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

15h33

* what if $X^T X$ is not invertible?

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

could choose $\hat{\omega} = \underset{\omega: X^T X \omega = X^T y}{\text{argmin}} \|\omega\| = X^+ y$ Moore-Penrose pseudo-inverse

problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect

regularization (motivated from MAP point of view)

$d \times d$ identity matrix

suppose we put prior $p(\omega) = N(\omega | 0, \frac{1}{\lambda^2} I)$

"precision" parameter

\leftarrow multivariate $N(\vec{\mu}, \Sigma_{d \times d})$

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

$$X^T y = \underbrace{\quad}_{\text{w.r. to } w} = -\frac{\|y - Xw\|^2}{2\sigma^2} + f(\sigma^2) - \frac{\lambda \|w\|^2}{2\sigma^2} + \text{const.}$$

MAP here

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

"ridge regression"

same as "regularized" ERM

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

is "strongly convex" in w

$f(\cdot)$ is λ -strongly convex

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

unique solution

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

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

no problem for $d > n$

one comment:

good practice to either standardize features i.e. make each feature zero mean and unit variance empirical

or normalize features \leftarrow make x_i unit norm $\|x_i\|_2 = 1$ or scale features to $[0, 1]$ or $[-1, 1]$

lehor

scale features to $[0,1]$ or $[-1,1]$

logistic regression

sets up binary classification $\mathcal{Y} = \{0,1\}$, $X \in \mathbb{R}^d$

generative

↳ motivation: suppose only assumption is \exists a pdf (densities) in \mathbb{R}^d

$$p(x | Y=1) \text{ \& } p(x | Y=0) \quad \text{"class conditionals"}$$

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

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

$$= \frac{1}{1 + \exp(-f(x))}$$

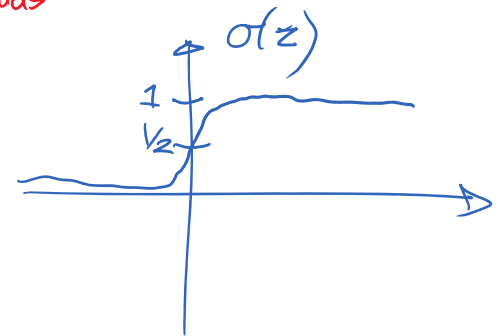
where $f(x) \triangleq \underbrace{\log \frac{P(X=x | Y=1)}{P(X=x | Y=0)}}_{\text{class conditional ratio}} + \underbrace{\log \frac{P(Y=1)}{P(Y=0)}}_{\text{prior odds ratio}}$

"log odds"

and thus in general, $P(Y=1 | X=x) = \sigma(f(x))$

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

"sigmoid function"



Some properties of $\sigma(z)$:

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

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

(*) to motivate linear logistic, consider class conditionals in the exponential family
 linear in η scalar functions log partition fct.

$$p(x|\eta) \triangleq \underbrace{h(x)}_{\text{"canonical parameter"}} \exp(\underbrace{\eta^T T(x)}_{\text{linear in } \eta}) - \underbrace{A(\eta)}_{\text{"sufficient statistics"}}$$

these specify the family

log odds $f(z) = \log \frac{p(z|\eta_1)}{p(z|\eta_0)} + \log \frac{\pi}{1-\pi}$

parameter for class 0

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

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

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

$$\phi(z) = \begin{pmatrix} T(z) \\ 1 \end{pmatrix}$$

"feature map"

→ get logistic regression model

→ gen logistic regression model

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

feature map

quiz

- **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](#))