

## Logistic regression (LR)

- linear decision rule with different  $w, b$  than LDA
- posterior probabilities (= soft classification)  $p(Y=k|X)$
- idea: - use the generative model of LDA (RHS of Bayes formula) and transform it into the corresponding posterior (LHS of Bayes)
  - train  $w$  such that the error rate of the maximum a-posteriori rule (MAP) is minimized
 
$$\hat{Y}_i = \arg \max_k p(Y=k|X_i) \text{ with } p(Y=k|X) = \frac{p(X|Y=k)p(Y=k)}{p(X)}$$

- simplification: centralize the data before training:

$$\text{mean } \mu = \frac{1}{N} \sum_{i=1}^N X_i, \text{ new features } \tilde{X}_i = X_i - \mu$$

[To avoid cluttered notation: do not write  $\tilde{X}_i$ , but  $X_i$  with implicit assumption that  $X_i$  are already centered]

useful, ~~also~~ because  $b=0$  for 2-class classification

$$\hat{Y}_i = \text{sign}(w X_i^T) \text{ if } X_i \text{ are centered}$$

- simplification: balanced classes for  $C=2$   $N_{-1} = N_{+1} = N/2$

$$\Rightarrow \frac{1}{N_{-1}} \sum_{i: Y_i = -1} X_i = \mu_{-1} = -\mu_1 = -\frac{1}{N_{+1}} \sum_{i: Y_i = 1} X_i$$

$$\mu = \frac{1}{N} \sum_i X_i \stackrel{!}{=} 0 = \frac{1}{2N_{-1}} \sum_{i: Y_i = -1} X_i + \frac{1}{2N_{+1}} \sum_{i: Y_i = 1} X_i = \frac{1}{2} \mu_{-1} + \frac{1}{2} \mu_1$$



- calculate the posterior for LDA when  $N_{-1} = N_{+1}$  and  $p(Y=-1) = p(Y=+1) = \frac{1}{2}$   
 [empirical priors  $\hat{p}(Y=y) = \frac{N_y}{N}$  will always fulfill this for balanced classes)]

$$p(Y=1 | x) = \frac{p(x | Y=1) p(Y=1)}{p(x)} = \frac{p(x | Y=1) p(Y=1)}{p(x | Y=-1) p(Y=-1) + p(x | Y=1) p(Y=1)}$$

$$p(Y=+1 | x) = \frac{1}{1 + \frac{p(x | Y=-1)}{p(x | Y=+1)}}$$

always true for  $C=2$   
 and then  $p(Y=-1) = p(Y=+1)$

insert  $p(x | Y=y)$  as multi-variate Gaussian  $N(\mu_y, \Sigma_w)$

$$\frac{p(x | Y=-1)}{p(x | Y=+1)} = \frac{\frac{1}{\sqrt{\det(2\pi \Sigma_w)}} \exp(-\frac{1}{2} (x - \mu_{-1})^T \Sigma_w^{-1} (x - \mu_{-1}))}{\frac{1}{\sqrt{\det(2\pi \Sigma_w)}} \exp(-\frac{1}{2} (x - \mu_{+1})^T \Sigma_w^{-1} (x - \mu_{+1}))}$$

$$(x - \mu_y)^T \Sigma_w^{-1} (x - \mu_y) = x^T \Sigma_w^{-1} x - 2 x^T \Sigma_w^{-1} \mu_y + \mu_y^T \Sigma_w^{-1} \mu_y$$

$$\frac{p(x | Y=-1)}{p(x | Y=+1)} = \exp\left(-\frac{1}{2} x^T \Sigma_w^{-1} x + x^T \Sigma_w^{-1} \mu_{-1} - \frac{1}{2} \mu_{-1}^T \Sigma_w^{-1} \mu_{-1} + \frac{1}{2} x^T \Sigma_w^{-1} x - x^T \Sigma_w^{-1} \mu_{+1} + \frac{1}{2} \mu_{+1}^T \Sigma_w^{-1} \mu_{+1}\right)$$

$$= \exp\left(-x^T \Sigma_w^{-1} (\mu_{+1} - \mu_{-1})^T + \frac{1}{2} (\mu_{+1}^T \Sigma_w^{-1} \mu_{+1} - \mu_{-1}^T \Sigma_w^{-1} \mu_{-1})\right)$$

$$= \exp(-x^T \Sigma_w^{-1} (\mu_{+1} - \mu_{-1})^T) \quad = 0 \text{ because } \mu_{+1} = -\mu_{-1}$$



$$p(Y=+1|x) = \frac{1}{1 + \frac{p(x|Y=-1)}{p(x|Y=+1)}} = \frac{1}{1 + \exp\left(-x \underbrace{\Sigma_w^{-1}(\mu_+ - \mu_-)^T}_{w^T}\right)}$$

$$w = (\mu_+ - \mu_-) \Sigma_w^{-1}$$

(remember:  $x_i$  are centered)

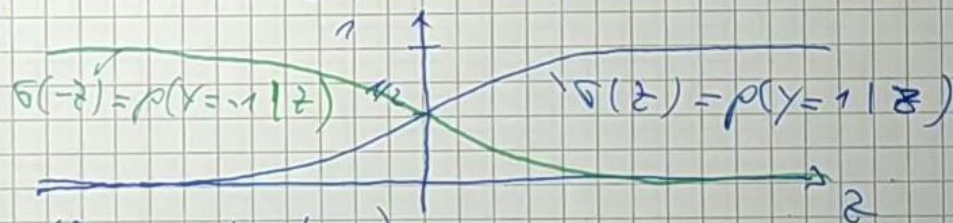
$$= \frac{1}{1 + \exp(-w^T x)} = p(Y=+1|x)$$

The function  $\sigma(z) = \frac{1}{1 + \exp(-z)}$  is called "logistic sigmoid function"

sigmoid  $\sigma$  graph looks like an "S"

logistic  $\hat{=}$  posterior of logistic regression

(other sigmoid exist, e.g. error function, tanh)



properties:  $p(Y=+1|x) = \sigma(w^T x)$

$$p(Y=-1|x) = 1 - p(Y=+1|x) = 1 - \sigma(w^T x) = \sigma(-w^T x)$$

e.g.  $z = w^T x \geq 0 \Rightarrow$  vote for class +1

$< 0 \Rightarrow$  vote for class -1

$$\text{[proof: } 1 - \frac{1}{1 + \exp(-z)} = \frac{1 + \exp(-z) - 1}{1 + \exp(-z)} = \frac{\exp(-z)}{1 + \exp(-z)} = \frac{1}{\exp(z) + 1} = \sigma(-z) \text{]}$$

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



Logistic regression continued

- assume centered data  $\frac{1}{N} \sum_{i=1}^N X_i = 0$ , balanced classes  $N_1 = N_0 = \frac{N}{2}$   $Y_i \in \{0, 1\}$
- Bayes formula:  $p(Y=k|X) = \frac{p(X|Y=k) p(Y=k)}{p(X)}$  with  $p(X|Y=k) = \mathcal{N}(\mu_k, \Sigma_w)$   
 common co-variance for all classes

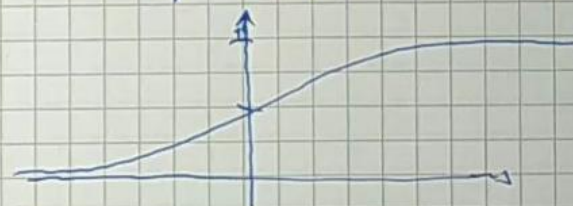
then the posterior becomes

$$p(Y=1|X_i) = \sigma(w X_i^T) = \frac{1}{1 + \exp(-w X_i^T)}$$

$$p(Y=0|X_i) = \sigma(-w X_i^T)$$

- LOA: choose  $w = (\mu_1 - \mu_0) \cdot \Sigma_w^{-1}$

result of optimal fit of RHS of Bayes formula



- Logistic regression (LR) optimally fits the LHS of Bayes

maximum likelihood principle:  $\hat{w} = \arg \max_w p(Y_1, \dots, Y_N | X_1, \dots, X_N)$

simplify via the i.i.d. assumption: joint probability factorizes

$$p(Y_1, \dots, Y_N | X_1, \dots, X_N) = \prod_{i: Y_i=0} p(Y=0 | X_i) \cdot \prod_{i: Y_i=1} p(Y=1 | X_i)$$

$$= \prod_{i=1}^N p(Y=0 | X_i)^{(1-Y_i)} \cdot p(Y=1 | X_i)^{Y_i}$$

trick: use  $Y_i \in \{0, 1\}$

take negative logarithm

$$-\log p(Y_1, \dots, Y_N | X_1, \dots, X_N) = - \sum_{i=1}^N \left[ (1-Y_i) \log \frac{p(Y=0 | X_i)}{\sigma(-w X_i^T)} + Y_i \log \frac{p(Y=1 | X_i)}{\sigma(w X_i^T)} \right]$$

optimization objective:  $\hat{w} = \arg \min_w \sum_{i=1}^N \left[ (1-Y_i) \sigma(-w X_i^T) + Y_i \sigma(w X_i^T) \right]$



optimization objective:  
of LP

$$\hat{w} = \arg \min_w - \sum_{i=1}^n \left[ (1-y_i) \log \sigma(-w x_i^T) + y_i \log \sigma(w x_i^T) \right]$$

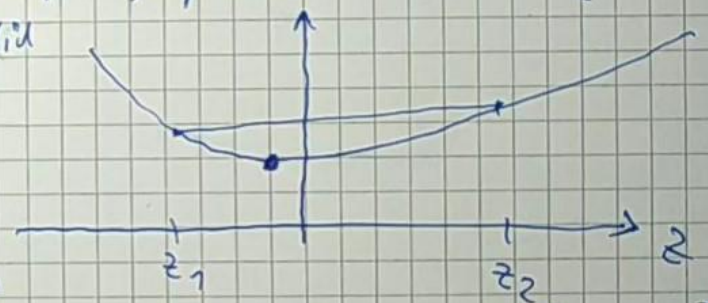
turns out that we now get a different  $\hat{w}$  than LDA

• but there is no closed-form expression for  $\hat{w} \Rightarrow$  need iterative alg.

iterative alg. will always find the optimal solution, because the objective is ~~convex~~

convex [  $f(z)$  is convex if  $f(\tau z_1 + (1-\tau)z_2) \leq \tau f(z_1) + (1-\tau)f(z_2)$   
with  $\tau \in [0,1]$ , intuition: the straight line between  $f(z_1)$  and  $f(z_2)$   
is always above the graph of  $f(z)$  between  $z_1$  and  $z_2$

strictly convex:  $\tau \in (0,1)$  the " $<$ " is valid  
 $\Rightarrow$  function  $f$  has only one minimum,  
which is automatically the global one



• to optimize, take derivative w.r.t.  $w$  and set to zero

derivative of logistic function  $\sigma'(z) = \frac{d}{dz} \frac{1}{1+e^{-z}} = \frac{1}{1+z} (1+e^{-z})^{-1} = \frac{1}{1+z} (1+e^{-z})^{-2} (-1)e^{-z}$

$$= \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{1}{1+e^z}$$

$$= \sigma(z) \cdot \sigma(-z)$$

$$\frac{d}{dw} (1-y_i) \log \sigma(-w x_i^T) = (1-y_i) \frac{1}{\sigma(-w x_i^T)} \cdot \frac{\sigma(-w x_i^T) \sigma(w x_i^T)}{\sigma(-w x_i^T)} (-x_i)$$

$$= \sigma'(-w x_i^T) (-x_i)$$



$$\frac{\partial}{\partial w} y_i \log \sigma(w x_i^T) = y_i \frac{1}{\sigma(w x_i^T)} \sigma(w x_i^T) \cdot \sigma(-w x_i^T) \cdot x_i$$

$$\begin{aligned} \frac{\partial}{\partial w} - \sum_{i=1}^N \dots &= - \sum_{i=1}^N \left[ (1-y_i) \sigma(w x_i^T) (-x_i) + y_i \frac{\sigma(-w x_i^T) x_i}{(1-\sigma(w x_i^T))} \right] \\ &= - \sum_{i=1}^N \sigma(w x_i^T) (-x_i) + y_i \sigma(w x_i^T) (+x_i) + y_i x_i - y_i \sigma(w x_i^T) x_i \\ &= \boxed{\sum_{i=1}^N (\sigma(w x_i^T) - y_i) x_i \stackrel{!}{=} 0} \quad \text{cannot analytically solve for } w \end{aligned}$$

i)  $\sigma(w x_i^T) \approx 1$  for instances with  $y_i = 1 \Rightarrow (\sigma(w x_i^T) - y_i) \approx 0 \Rightarrow$  gradient vanishes  
 $\hat{=} w x_i^T \rightarrow \infty$

ii)  $\sigma(w x_i^T) \approx 0$  for instances with  $y_i = 0 \Rightarrow (\sigma(w x_i^T) - y_i) \approx 0 \Rightarrow$  gradient vanishes  
 $\hat{=} w x_i^T \rightarrow -\infty$

if an  $y_i = 1$  instance is wrongly classified by the current guess  $w \hat{=} \sigma(w x_i^T) < 1$   
 $\Rightarrow$  gradient  $(\sigma(w x_i^T) - y_i)$  is negative  $\Rightarrow$  change  $w x_i^T$  in the opposite direction to make it bigger  $\Rightarrow \sigma(w x_i^T)$  gets closer to 1

if an  $y_i = 0$  instance is wrongly classified  $\Rightarrow (\sigma(w x_i^T) - y_i) > 0 \Rightarrow$  pulls in the opposite direction to make  $\sigma(w x_i^T)$  smaller

$\Rightarrow$  solve by gradient descent • initial guess  $w_0 = 0$ , learning rate  $\eta$   
 • for  $t = 1, \dots, T$  # of iterations:  $w_t = w_{t-1} - \eta \sum_{i=1}^N (\sigma(w_{t-1} x_i^T) - y_i) x_i$   
 minus, because we want to minimise the objective  $\Rightarrow$  gradient descent



- variations of gradient descent:
- mini-batch descent: only include a ~~finite~~ random subset of TS in the sum  
(choose  $B$  instances uniformly at random  $\Rightarrow$  mini-batch)  
faster per iteration (only  $O(B)$  instead  $O(N)$ , but more iterations)
  - stochastic gradient descent:  $B=1$ , i.e. approximate sum by a single random instance (SGD)  
 $\Rightarrow$  very fast per iteration, but even more iterations
  - learning rate schedule: decrease  $\rightarrow$  regularity, e.g.  $\eta = \frac{\eta_0}{t+c}$  -  $\eta_0$   $\leftarrow$  initial  $\eta$   
(especially important for mini-batch SGD for convex costs)  $\leftarrow$  iteration counter
  - in practice: split loop into two
 

for  $e = 1, \dots, E$        $e$ : epoch       $E$  number of epochs  
 shuffle the TS  
 for  $i = 1, \dots, N$ :  
     update  $w$  with the gradient of instance  $\tilde{\pi}(i)$ 

permutation after shuffle  
 $\downarrow$
- faster algorithm: Newton alg. or quasi-Newton alg.  
needs much fewer iterations, ~~but~~ but each iteration is expensive
  - fastest alg. depends on size of training set  $N$ 
    - $N$  small: Newton is fastest (few iterations, relatively cheap iterations)
    - $N$  large: (mini-batch) SGD is fastest (many iterations, but compensated because ~~each~~ each iteration is very cheap)