

# ML solution

## Lecture 8

Recall that we want to use

$$\mathcal{D} = \{\vec{x}_n, \vec{t}_n\}$$

dataset  
 $n=1, \dots, N$

$$p(c_i | \vec{x}, \mathcal{D}) = \frac{p(\vec{x} | c_i, \mathcal{D}) p(c_i | \mathcal{D})}{\sum_{j=1}^K p(\vec{x} | c_j, \mathcal{D}) p(c_j | \mathcal{D})}$$

Before we assumed that  $p(\vec{x} | c_i)$  is a known gaussian and that  $p(c_i)$  is also given. now let's estimate them by ML.

so formally speaking they do not depend on  $\mathcal{D}$

Specifically, consider  $K=2$  for simplicity.

$$\text{Define } \begin{cases} p(c_1 | \mathcal{D}) = \pi & \Rightarrow p(c_2 | \mathcal{D}) = 1 - \pi; \\ t_n = 1 \text{ in } c_1, t_n = 0 \text{ in } c_2. \end{cases}$$

Then  $\mathcal{Z} = \prod_{n=1}^N \pi^{t_n} (1-\pi)^{1-t_n}$ , or

$$\log \mathcal{Z} = \sum_{n=1}^N \left\{ t_n \log \pi + (1-t_n) \log(1-\pi) \right\}.$$

$$\frac{\partial}{\partial \pi} \log \mathcal{Z} = \frac{1}{\pi} \sum_{n=1}^N t_n - \frac{1}{1-\pi} \sum_{n=1}^N (1-t_n) = 0, \text{ or}$$

$$(1-\pi) \sum_{n=1}^N t_n - \pi \sum_{n=1}^N (1-t_n) = 0,$$

$$\pi N = \sum_{n=1}^N t_n \Rightarrow \pi = \frac{1}{N} \sum_{n=1}^N t_n$$

"  $p(c_1 | \mathcal{D})_{ML}$

Likewise, assume

$$\begin{cases} p(\vec{x}_n | C_1) = \mathcal{N}(\vec{x}_n | \vec{\mu}_1, \Sigma) \\ p(\vec{x}_n | C_2) = \mathcal{N}(\vec{x}_n | \vec{\mu}_2, \Sigma) \end{cases}$$

← shared

estimate  $\vec{\mu}_1, \vec{\mu}_2, \Sigma$  by ML

Then

$$\mathcal{J}' = \prod_{n=1}^N [\mathcal{N}(\vec{x}_n | \vec{\mu}_1, \Sigma)]^{t_n} [\mathcal{N}(\vec{x}_n | \vec{\mu}_2, \Sigma)]^{1-t_n},$$

$$\begin{aligned} \log \mathcal{J}' &= -\frac{1}{2} \sum_{n=1}^N t_n (\vec{x}_n - \vec{\mu}_1)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_1) - \\ &\quad - \frac{1}{2} \sum_{n=1}^N (1-t_n) (\vec{x}_n - \vec{\mu}_2)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_2) - \\ &\quad - \frac{1}{2} \sum_{n=1}^N t_n \log |\Sigma| - \frac{1}{2} \sum_{n=1}^N (1-t_n) \log |\Sigma|. \end{aligned}$$

$$\frac{\partial}{\partial \vec{\mu}_1} \log \mathcal{J}' = 0 \Rightarrow \sum_n t_n \Sigma^{-1} (\vec{x}_n - \vec{\mu}_1) = 0,$$

$$\Sigma^{-1} \left( \sum_n t_n (\vec{x}_n - \vec{\mu}_1) \right) = 0$$

$$\vec{\mu}_1 = \frac{\sum_n t_n \vec{x}_n}{\underbrace{\sum_n t_n}_{N_1}} = \frac{1}{N_1} \sum_n t_n \vec{x}_n$$

Likewise,  $\vec{\mu}_2 = \frac{1}{N_2} \sum_n (1-t_n) \vec{x}_n$ .

Finally,

$$\log \mathcal{J}' = -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{n \in C_1} (\vec{x}_n - \vec{\mu}_1)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_1) - \frac{1}{2} \sum_{n \in C_2} (\vec{x}_n - \vec{\mu}_2)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_2) + \text{const}(\Sigma) \quad \textcircled{11}$$

$$\begin{aligned} \sum_{n \in C_1} \sum_{i,j} (x_{n,i} - \mu_{1,i}) \Sigma^{-1}_{ij} (x_{n,j} - \mu_{1,j}) &= \\ = \sum_{i,j} \Sigma^{-1}_{ij} \left[ \sum_{n \in C_1} (x_{n,j} - \mu_{1,j}) (x_{n,i} - \mu_{1,i}) \right] &= \\ = \text{Tr}(\Sigma^{-1} S_1) \cdot N_1. \end{aligned}$$

$$\begin{aligned} \textcircled{12} \quad -\frac{N}{2} \log |\Sigma| - \frac{N}{2} \text{Tr}(\Sigma^{-1} \underbrace{[\frac{N_1}{N} S_1 + \frac{N_2}{N} S_2]}_S) &= \\ = -\frac{N}{2} \log |\Sigma| - \frac{N}{2} \text{Tr}(\Sigma^{-1} S). \end{aligned}$$

$$\frac{\partial}{\partial \Sigma} \log \mathcal{J}' = -\frac{N}{2} \Sigma^{-1} + \frac{N}{2} (\Sigma^{-1} S \Sigma^{-1}) \stackrel{!}{=} 0, \text{ so that } \Sigma^{-1} = \Sigma^{-1} S \Sigma^{-1} \Rightarrow \boxed{S = \Sigma}$$

$$\stackrel{(C.28)}{\nearrow} \frac{\partial}{\partial \Sigma} \log |\Sigma| = \underbrace{\Sigma^{-1}}_{(\Sigma^{-1})^T}$$

$$\begin{aligned} \frac{\partial}{\partial \Sigma_{ij}} \text{Tr}(\Sigma^{-1} S) &= \text{Tr} \left( \frac{\partial}{\partial \Sigma_{ij}} (\Sigma^{-1}) S \right) = \\ = \text{Tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \Sigma_{ij}} \Sigma^{-1} S \right) &= -\text{Tr} \left( \mathbb{I}_{ij} \Sigma^{-1} S \Sigma^{-1} \right) = \\ = -(\Sigma^{-1} S \Sigma^{-1})_{ij}. \end{aligned}$$

1 at (i,j),  
0 everywhere else  
[ignores  $\Sigma_{ij} = \Sigma_{ji}$ ]

Finally,  $p(c_1|\vec{x}) = \sigma(\vec{w}^T \vec{x} + w_0)$ , with

$\vec{w}$  &  $w_0$  f's of  $\vec{\mu}_1, \vec{\mu}_2, \Sigma$  which are all estimated by ML.

$$p(c_2|\vec{x}) = 1 - p(c_1|\vec{x}).$$

Alternatively, we do not have to use the Bayes formula and can find  $p(c_k|\vec{x})$  more directly. This is often easier and may lead to better fits.

We will use a fixed non-linear transform into feature space:  $\vec{y}(\vec{x})$ .

### Logistic regression

Consider  $(K=2)$  again:

$$\begin{cases} p(c_1|\vec{y}) = y(\vec{y}) = \sigma(\vec{w}^T \vec{y}) \\ p(c_2|\vec{y}) = 1 - p(c_1|\vec{y}) \end{cases} \quad \uparrow y_0(\vec{x}) = 1$$

If  $\vec{y}$  has dimensions  $M$ , there are  $M$  fitting prms. Recall that in the previous approach, we need  $2M$  prms to fit  $\vec{\mu}_1, \vec{\mu}_2$ ,  $\frac{M(M+1)}{2}$  prms to fit  $\Sigma$  (shared cov. matrix), and 1 prm to fit  $P(c_1)/P(c_2)$ . So it's  $\Theta(M)$  vs.  $\Theta(M^2)$ .

Consider  $\{\vec{y}_n, t_n\} \quad n=1, \dots, N$   
 $t_n \in \{0, 1\} \quad \vec{y}_n = \vec{g}(\vec{x}_n)$

Then  $p(\vec{t} | \vec{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$

$\vec{t} \quad \left\{ \begin{array}{l} y_n = p(c_1 | \vec{y}_n), \\ 1-y_n = p(c_2 | \vec{y}_n). \end{array} \right. = \sigma(\underbrace{\vec{w}^T \vec{y}_n}_{a_n})$

The error f'n is given by


$$E(\vec{w}) = -\log p(\vec{t} | \vec{w}) = -\sum_{n=1}^N [t_n \log y_n + (1-t_n) \log(1-y_n)]$$

$$\frac{\partial E(\vec{w})}{\partial w_i} = -\sum_{n=1}^N \left[ \frac{t_n}{y_n} \frac{\partial y_n}{\partial w_i} - (1-t_n) \frac{1}{1-y_n} \frac{\partial y_n}{\partial w_i} \right] =$$

$$= \sum_{n=1}^N \underbrace{\frac{(1-t_n)y_n - t_n(1-y_n)}{y_n(1-y_n)}}_{\frac{y_n - t_n}{y_n(1-y_n)}} \underbrace{\frac{\partial y_n}{\partial w_i}}_{y_n(1-y_n) \frac{\partial a_n}{\partial w_i}} \quad \textcircled{=}$$

$$\textcircled{=} \sum_{n=1}^N (y_n - t_n) \underline{\underline{g_{n,i}}}$$

Difference between  $t_n$  &  $y_n$  multiplied by the basis f'n values.

Note that  $E(\vec{w})$  is a non-linear f'n of  $\vec{w}$  so there's no closed-form solution (unlike regression). However, one can argue that  $E(\vec{w})$  is ~~convex~~<sup>convex</sup> everywhere:  has a unique  $\vec{w}$  minimum.

We can find it using the Newton-Raphson (NR) algorithm:

$$\vec{\nabla} \equiv \frac{\partial}{\partial \vec{w}}$$

$$\vec{w}^{\text{new}} = \vec{w}^{\text{old}} - H^{-1} \vec{\nabla} E(\vec{w})$$

$H$  is the Hessian matrix:  $H_{ij} = \frac{\partial^2 E(\vec{w})}{\partial w_i \partial w_j}$

Let's first apply it to regression:

$$\tilde{E}(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{y}_n)^2 \quad \text{quadratic error f'n}$$

$$\text{Then } \frac{\partial \tilde{E}}{\partial w_i} = \sum_{n=1}^N (\vec{w}^T \vec{y}_n - t_n) y_{n,i} \Rightarrow$$

$$\Rightarrow \frac{\partial \tilde{E}}{\partial \vec{w}} = (\Phi^T \Phi) \vec{w} - \Phi^T \vec{t}$$

Further,

$$\frac{\partial^2 \tilde{E}}{\partial w_i \partial w_j} = \sum_{n=1}^N y_{n,i} y_{n,j} \Rightarrow H = \Phi^T \Phi$$

$$\begin{aligned} \text{Then } \vec{w}^{\text{new}} &= \vec{w}^{\text{old}} - (\Phi^T \Phi)^{-1} [(\Phi^T \Phi) \vec{w}^{\text{old}} - \Phi^T \vec{t}] = \\ &= \underline{\underline{(\Phi^T \Phi)^{-1} \Phi^T \vec{t}}} \quad \text{exact sol'n in one step} \end{aligned}$$

NR gives the exact solution b/c  $\tilde{E}(\vec{w})$  is quadratic.

Now let's consider  $E(\vec{w})$ :

$$\vec{\nabla} E = \Phi^T (\vec{y} - \vec{t}) \leftarrow \begin{array}{l} \text{non-linear} \\ \text{function of } \vec{w} \end{array}$$

$$\vec{y} = \{y_1, \dots, y_N\} \quad \vec{t} = \{t_1, \dots, t_N\}$$

$$\vec{H} = \vec{\nabla} \vec{\nabla} E \Rightarrow \frac{\partial E}{\partial w_i \partial w_j} = \sum_{n=1}^N \frac{\partial y_n}{\partial w_j} y_{n,i} =$$

$$= \sum_n y_n (1 - y_n) y_{n,j} y_{n,i}$$

Then

$$H = \Phi^T R \Phi$$

$$R = \begin{pmatrix} y_1(1-y_1) & & & \\ & \ddots & & \\ & & y_N(1-y_N) & \\ & & & \ddots \end{pmatrix} \left. \vphantom{\begin{pmatrix} \\ \\ \\ \end{pmatrix}} \right\} N$$

diagonal matrix

Note that  $H = H(\vec{w})$  through  $y_n$ .

since  $E(\vec{w})$  is non-quadratic

Since  $y_n(1-y_n) > 0, \forall n \Rightarrow H$  is positive definite.

But then  $E(\vec{w})$  is convex everywhere & has a unique minimum.

Finally, NR gives

$$\begin{aligned} \vec{w}^{\text{new}} &= \vec{w}^{\text{old}} - (\Phi^T R \Phi)^{-1} \Phi^T (\vec{y} - \vec{t}) = \\ &= (\Phi^T R \Phi)^{-1} [(\Phi^T R \Phi) \vec{w}^{\text{old}} - \Phi^T (\vec{y} - \vec{t})] = \end{aligned}$$

$$= (\Phi^T R \Phi)^{-1} \Phi^T R \vec{z} \leftarrow \begin{array}{l} N\text{-dim. vector} \\ \vec{z} = \vec{w}^{\text{old}} - R^{-1}(\vec{y} - \vec{t}) \end{array} \quad (*)$$

Apply (\*) iteratively to get  $\vec{w}$ .

Note that, as with regression,

$$E[t|\vec{x}] = y(\vec{x}) = \sigma(\underbrace{\vec{w}^T \cdot \vec{y}(\vec{x})}_{\text{fitted weights}})$$

$$\text{Var}[t|\vec{x}] = E[t^2|\vec{x}] - E^2[t|\vec{x}] = y - y^2 = y(1-y).$$

"   
 t since  $t \in \{0,1\}$

So, elements of the R matrix are variances of  $t_n$ .

~~Moreover, note that for  $d_n = \frac{1}{n} \sum_{i=1}^n t_i$~~   
 ~~$d_n(t)$~~



Alternatively, consider

$$\vec{\nabla} E(\vec{w}) = \Phi^T (\vec{y} - \vec{t}) = 0 \quad \text{directly.}$$

# prms  $\rightarrow M \times N \leftarrow$  # datapoints  
 $\uparrow \quad \uparrow$   
 rows columns

note that here  $\Phi^T \vec{y} = \Phi^T \vec{t}$  leads to

$$\vec{y} = (\Phi^T)^+ \Phi^T \vec{t}$$

Moore-Penrose pseudoinverse  
 $(\Phi \Phi^T)^{-1} \Phi$ , same as below

Then  $(\Phi \Phi^T) \vec{y} = (\Phi \Phi^T) \vec{t}$ , or

$$[N \times M * M \times N] = N \times N$$

$\equiv \mathbb{I}$

$$\vec{y} = (\Phi \Phi^T)^{-1} (\Phi \Phi^T) \vec{t} \Rightarrow \vec{y} = \vec{t},$$

as expected

$$\{ \sigma(a_1) \dots \sigma(a_N) \}$$

In components,

$$y_n = \sigma(a_n) = \frac{t_n}{1 + e^{-a_n}}$$

but then

$$a_n = \log \left( \frac{t_n}{1 - t_n} \right) \equiv C_n$$

component of an N vector with  $\{+\infty, -\infty\}$  components (!)

$$\sum_{j=1}^M w_j y_{n,j} = \sum_j \phi_{nj} w_j$$

$y_j(\vec{x}_n) = \phi_{nj}$

Finally, in matrix form

$$\Phi \vec{w} = \vec{C} \Rightarrow \vec{w} = \Phi^+ \vec{C} = (\Phi^T \Phi)^{-1} \Phi^T \vec{C}$$

$\Phi$ :  $M \times N$  Mvec  
 $\vec{w}$ :  $N$  vec  
 $\vec{C}$ :  $N$  vec  
 $(\Phi^T \Phi) \vec{w}$ :  $N$  vector

Clearly at least some weights will be  $\pm\infty$ , but the solution is unique