

Chapter 3Linear models for classificationI) Logistic regression

Remember that the model is the following:

$$p(y|x, w) = \mathcal{B}(y | \text{sigma}(w^T x)) \quad \text{and, sometimes, this is replaced by } \text{sigma}(w_0 + w^T x)$$

$\rightarrow y$ is a Bernoulli law of parameter $\text{sigma}(w^T x)$

1) MLE

The negative log-likelihood for logistic regression is given by:

$$\begin{aligned} \text{NLL}(w) &= - \sum_{i=1}^N \log(p_i^{y_i=1} (1-p_i)^{1-y_i=0}) \\ &\quad \rightarrow \text{where } p_i = \text{sigma}(w^T x_i) \\ &= - \sum_{i=1}^N [y_i \log p_i + (1-y_i) \log(1-p_i)] \\ &\quad (\text{remember: } \text{sigma}(y) = \frac{1}{1+e^{-y}}) \end{aligned}$$

Let us set $\tilde{y}_i = 1$ if $y_i = 1$
 $\tilde{y}_i = -1$ if $y_i = 0$.

$$\text{We then have: } \begin{cases} \mathbb{P}(\tilde{y}_i = 1) = \frac{1}{1 + \exp(-w^T x_i)} \\ \mathbb{P}(\tilde{y}_i = -1) = \frac{1}{1 + \exp(w^T x_i)} \end{cases}$$

$$\text{and so: } \text{NLL}(w) = \sum_{i=1}^N \log(1 + \exp(-\tilde{y}_i w^T x_i)).$$

There is no closed form expression for the MLE. We need to use an optimization algorithm to compute it.

Let us set $\sigma(\dots) = \text{sign}(\dots)$.

We have $\therefore \sigma'(t) = \frac{e^{-t}}{(1+e^{-t})^2} = \frac{1}{1+e^{-t}} \left(1 - \frac{1}{1+e^{-t}}\right)$
 $= \sigma(t)(1-\sigma(t))$

So :

$$\begin{aligned} \underbrace{\nabla NLL(w)}_{=: g(w)} &= - \sum_{i=1}^N \left\{ \frac{y_i}{p_i} \sigma'(w^T x_i) x_i - \frac{(1-y_i)}{1-p_i} \sigma'(w^T x_i) x_i \right\} \\ &= - \sum_{i=1}^N \left\{ \left(\frac{y_i}{p_i} p_i (1-p_i) x_i - \frac{(1-y_i)}{1-p_i} p_i (1-p_i) x_i \right) \right\} \\ &= - \sum_{i=1}^N \{ (y_i - p_i) x_i \} \end{aligned}$$

$$g(w) = \sum_{i=1}^N (p_i - y_i) x_i$$

We write : $x_i = \begin{pmatrix} x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,D} \end{pmatrix}$ $X = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ \vdots & \vdots & & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix}$

(the space is \mathbb{R}^D and we have N observations)

$$\text{So } \nabla NLL(w) = \underbrace{\begin{bmatrix} x_{1,1} & \dots & x_{N,1} \\ x_{1,2} & & x_{N,2} \\ \vdots & & \vdots \\ x_{1,D} & & x_{N,D} \end{bmatrix}}_{= X^T} \times \underbrace{\begin{pmatrix} p_1 - y_1 \\ \vdots \\ p_N - y_N \end{pmatrix}}_{= p - y}$$

$$H(w) = \sum_{i=1}^N (D_w p_i) x_i^T = \sum_{i=1}^N p_i (1-p_i) x_i x_i^T = X^T S X$$

where $S := \text{diag}(p_i (1-p_i))_{1 \leq i \leq N}$

Suppose that $\forall i, p_i \in (0,1)$ and that $N > D$,
 $\text{rank}(X) = D$. In this case:

29

We write $\tilde{x}_j = \begin{pmatrix} x_{1,j} \\ x_{2,j} \\ \vdots \\ x_{N,j} \end{pmatrix}$ for all $j \in \{1, 2, \dots, D\}$.

For any $g = \begin{pmatrix} g_1 \\ \vdots \\ g_D \end{pmatrix} \in \mathbb{R}^D$, we have: $Xg = [\tilde{x}_1 | \dots | \tilde{x}_D] \begin{pmatrix} g_1 \\ \vdots \\ g_D \end{pmatrix}$
 $= g_1 \tilde{x}_1 + \dots + g_D \tilde{x}_D$.

So: $g^T X^T S X g = (g_1 \tilde{x}_1 + \dots + g_D \tilde{x}_D)^T S (g_1 \tilde{x}_1 + \dots + g_D \tilde{x}_D)$
 $= \sum_{i=1}^N p_i (1-p_i) \left(\sum_{j=1}^D g_j x_{i,j} \right)^2$

↳ It is always ≥ 0 .

And if it is equal to zero then:

$$\forall i \in \{1, 2, \dots, N\} : \sum_{j=1}^D g_j x_{i,j} = 0$$

$$\text{that is: } \sum_{j=1}^D g_j \tilde{x}_j = 0$$

(as $\tilde{x}_1, \dots, \tilde{x}_D$ are linearly independent)

this proves $g_1 = g_2 = \dots = g_D = 0$.

So H is definite positive.

2) Optimization algorithms (non exhaustive list)

a) Steepest descent (/ gradient descent)

We define a sequence:

$$\begin{cases} \theta_0 \in \mathbb{R}^D \\ \theta_{k+1} = \theta_k - \eta_k g(\theta_k) \end{cases}$$

step size / learning rate

How should we set the step size?

- If it is too small, the convergence will be slow.
- If it is too big, it might never converge.

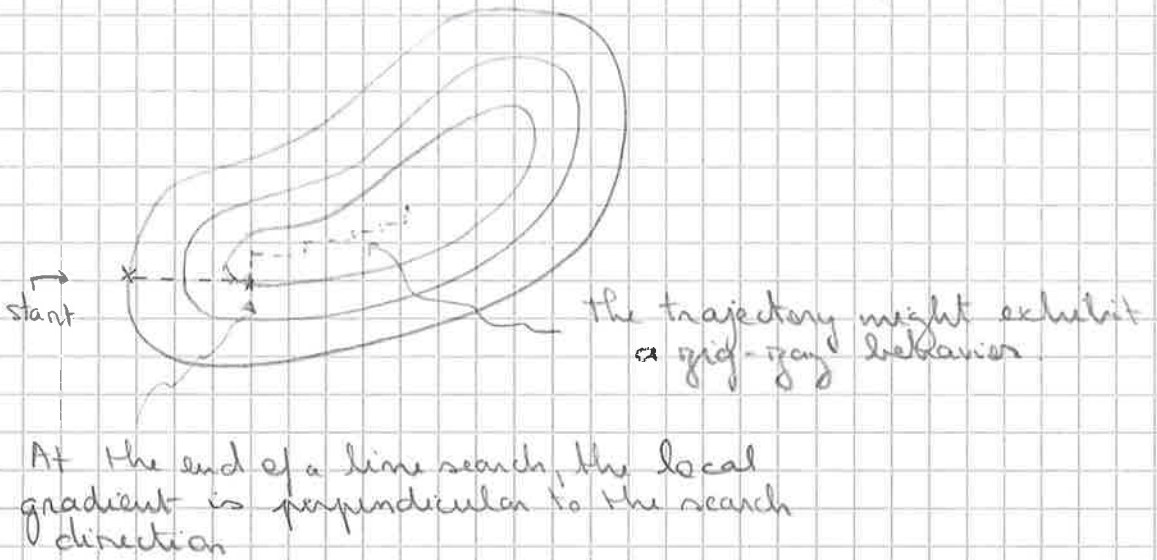
We observe that: $f(\theta + \eta d) \approx f(\theta) + \eta \nabla f(\theta) d$
 function we want to minimize Taylor expansion (where d is an descent direction)

If η is small enough then $f(\theta + \eta d) < f(\theta)$.

Let us pick η to minimize: $\phi(\eta) = f(\theta + \eta d)$.

(This is called line minimization / line search.)

▷ This is a 1-d minimization problem which can be solved using various methods.



1) Newton's method (for a strictly convex function)

We could obtain a faster optimization method by taking into account the curvature of the space (i.e. the Hessian).

These are called second order optimization methods. The primary example is Newton's algorithm.

We set :

$$\begin{cases} \theta_0 \in \mathbb{R}^d \\ \theta_{k+1} = \theta_k - \eta_k H_k^{-1} g_k \end{cases}$$

where :

$$\begin{cases} H_k = H(\theta_k) \\ g_k = \nabla f(\theta_k) \end{cases}$$

(H is the Hessian of f ,
 f is the function to minimize)

$$\begin{cases} d_k \text{ sol. of } H_k d_k = -g_k \\ \eta_k = \underset{\eta}{\operatorname{argmin}} f(\theta_k + \eta d_k) \end{cases}$$

the idea behind the algorithm is the following.

$$f(\theta) \underset{\uparrow}{=} f(\theta_k) + g_k^T (\theta - \theta_k) + \frac{1}{2} (\theta - \theta_k)^T H_k (\theta - \theta_k) + \dots$$

second order Taylor expansion

$$f(\theta) = c + b^T \theta + \theta^T A \theta$$

\downarrow $f_k - g_k^T \theta_k$ \downarrow $b = g_k - H_k \theta_k$ \downarrow $A = \frac{1}{2} H_k$

The min. of this is at : $\theta = -\frac{1}{2} A^{-1} b = \theta_k - H_k^{-1} g_k$

thus the Newton step $d_k = -H_k^{-1} g_k$ is what should be added to θ_k to minimize the second order approximation of f around θ_k . We choose then η_k (not necessarily = 1) using a line search to be even more optimal.

Application of Newton's method : IRLS

We want to apply the above method to find the MLE for binary logistic regression. We fix $\eta_k = 1$ (as the Hessian is constant, the Taylor expansion of order 2 is exact, so the best θ is given for $\eta_k = 1$).

We have:

$$\theta_{k+1} = \theta_k - H^{-1} g_k$$

27

$$= \theta_k + (X^T S_k X)^{-1} X^T (y - p_k)$$

↪ at iteration k

$$= (X^T S_k X)^{-1} [(X^T S_k X) \theta_k + X^T (y - p_k)]$$

$$= (X^T S_k X)^{-1} X^T [S_k X \theta_k + y - p_k]$$

$$= (X^T S_k X)^{-1} X^T S_k \underbrace{z_k}_{z_k = X \theta_k + S_k^{-1} (y - p_k)}$$

called working response

Observe that θ_{k+1} is a minimizer of:

$$(D_k g_k - D_k X \theta)^T (D_k g_k - D_k X \theta)$$

where $S_k = D_k^2$

and: $(D_k g_k - D_k X \theta)^T (D_k g_k - D_k X \theta)$

$$\sum_{i=1}^N \sqrt{p_{ki}(1-p_{ki})} (r_{k,i} - \theta^T x_i)^2$$

which is called a weighted least square problem.

We can write for each component of g_k :

$$r_{k,i} = \theta_k^T x_i + \frac{y_i - p_{ki}}{p_{ki}(1-p_{ki})}$$

This algorithm is known as iterated reweighted least square

(IRLS) since at each iteration, we solve a weighted least squares problem, where the weight matrix S_k changes at each iteration.

↓
because
 $(X^T S_k X)^{-1} X^T S_k g_k$
"
 $(X^T D_k^T D_k X)^{-1} X^T D_k^T D_k g_k$
+ remember
chapter 2

Algorithm

$$\theta = 0 \quad ; \quad w_0 = \log(\bar{y} / (1 - \bar{y}))$$

repeat:

$$\left. \begin{aligned} \eta_i &= w_0 + \theta^T x_i \\ \mu_i &= \text{sign}(\eta_i) \\ s_i &= \mu_i (1 - \mu_i) \\ \tilde{y}_i &= \eta_i + \frac{y_i - \mu_i}{s_i} \end{aligned} \right\} \text{for all } i$$

$$S = \text{diag}(s_{1:n})$$

$$\theta = (X^T S X)^{-1} X^T S \tilde{y}$$

until converged

3) l_2 -regularization

If the data is linearly separable, the MLE is obtained when $\|w\| \rightarrow +\infty$, corresponding to an infinitely steep sigmoid function $\mathbb{1}_{w^T x > w_0}$ also known as linear threshold limit. This assigns the maximal amount of probability mass to the training data (which will not be good for generalisation).

To prevent this, we can use a l_2 -regularisation (just like in the ridge regression). The new objective function, gradient and Hessian have the form:

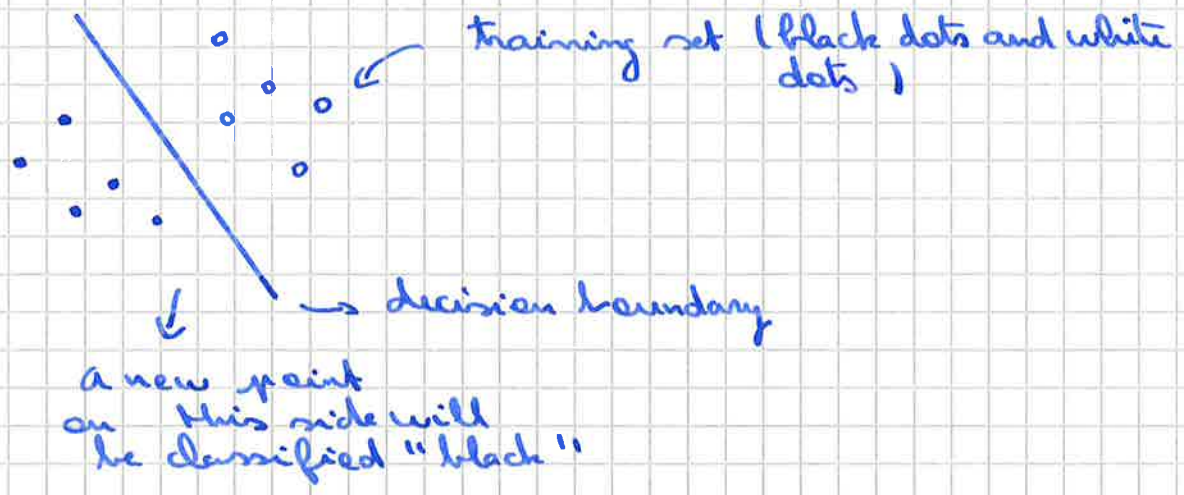
$$\begin{cases} \tilde{J}(w) = \text{NLL}(w) + \lambda w^T w \\ \tilde{g}(w) = g(w) + \lambda w \\ \tilde{H}(w) = H(w) + \lambda \text{Id} \end{cases}$$

(again, we can apply our optimization alg. to \tilde{J}).

II) Linear support vector machine (SVM)

The linear SVM uses an hyperplane as a decision boundary

Ex:



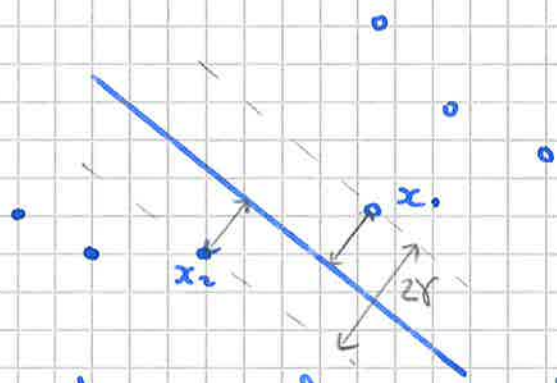
The decision boundary is $\{x : w^T x + b = 0\}$ for a vector w and a scalar b .

The SVM decision function for a test point x_{new} is therefore given by: $t_{new} = \text{sign}(w^T x_{new} + b)$
($t_{new} \in \{-1, +1\}$)

1) The margin

The margin is defined as the perpendicular distance from the decision boundary to the closest point on either side.

Ex: margin = γ



x_1, x_2 and the boundary are always such that $\text{dist}(x_1, \text{boundary}) = \text{dist}(x_2, \text{boundary})$

2) Maximizing the margin

w is a vector perpendicular to the boundary so

$$2\gamma = \frac{1}{\|w\|} w^T (x_1 - x_2)$$

The decision function $t_{\text{new}} = \text{sign}(w^T x_{\text{new}} + b)$ is invariant to the scaling of its argument by a positive constant. This means that we can multiply $w^T x_{\text{new}} + b$ by a positive λ and the output of the function will be unchanged.

Therefore we can fix the scaling of w and b such that $w^T x + b = \pm 1$ for the closest points on either side. Thus:

$$\begin{aligned} 2\gamma &= \frac{w^T}{\|w\|} (x_1 - x_2) = \frac{1}{\|w\|} (w^T x_1 + b - w^T x_2 - b) \\ &= \frac{(1+1)}{\|w\|} \end{aligned}$$

$$\gamma = \frac{1}{\|w\|}$$

To maximize the margin, we must therefore maximize $\frac{1}{\|w\|}$.

There are some constraints: $\left\{ \begin{array}{l} w^T x_n + b \geq 1 \text{ for all points in class 1} \\ w^T x_n + b \leq -1 \text{ for all points in class 2} \end{array} \right.$

We formalize the problem in the following form:

$$\textcircled{1} \left\{ \begin{array}{l} \text{find } \hat{w} = \underset{w}{\text{argmin}} \frac{1}{2} \|w\|^2 \\ \text{subject to } t_n (w^T x_n + b) \geq 1 \text{ for all } n \end{array} \right.$$

To solve this optimization problem with constraints, we need to incorporate the constraints into the objective function through a set of Lagrange multipliers.

Remark: See the optimization course (M1 in),
 for the Lagrange multipliers or
 [BD2001] Binnore, Ken and Davies, Jean; Calculus,
 concepts and methods

Our new objective function is:

$$\textcircled{2} \begin{cases} (w, \alpha) = \underset{w, \alpha}{\operatorname{argmin}} \frac{1}{2} w^T w - \sum_{n=1}^N \alpha_n (t_n (w^T x_n + b) - 1) \\ \text{subject to } \alpha_n \geq 0 \text{ for all } n \end{cases} = f_L(w, \alpha)$$

here, we have N
 points in the learning
 data set

(The idea of the Lagrange multiplier technique is that
 a sol. to problem ① is a solution to problem ②
 - and that we apply the reverse implication without
 proof.)

We look for a critical point of our new objective function:

$$\begin{cases} \nabla_w f_L = w - \sum_{n=1}^N \alpha_n t_n x_n \\ \frac{\partial f_L}{\partial b} = - \sum_{n=1}^N \alpha_n t_n \end{cases}$$

At a critical point:

$$\begin{cases} w = \sum_{n=1}^N \alpha_n t_n x_n \\ \sum_{n=1}^N \alpha_n t_n = 0 \end{cases}$$

If we compute the value of f_L at a critical point, we get:

$$\frac{1}{2} w^T w - \sum_{n=1}^N \alpha_n (t_n (w^T x_n + b) - 1) =$$

$$\frac{1}{2} \left(\sum_{n=1}^N \alpha_n t_n x_n \right)^T \left(\sum_{n=1}^N \alpha_n t_n x_n \right) - \sum_{n=1}^N \alpha_n \left(t_n \left(\sum_{m=1}^N \alpha_m t_m x_m^T x_n + b \right) - 1 \right) = \dots$$

$$\dots = \frac{1}{2} \sum_{n=1}^N \alpha_n \alpha_n t_n t_n x_n^T x_n - \sum_{n=1}^N \alpha_n \alpha_n t_n t_n x_n^T x_n - \sum_{n=1}^N \alpha_n t_n b + \sum_{n=1}^N \alpha_n$$

(as $\sum_{n=1}^N \alpha_n t_n = 0$)

$$= \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \alpha_n \alpha_n t_n t_n x_n^T x_n$$

There is a continuum of critical points. We want to find, among the critical points, the one minimizing the above expression:

$$\begin{cases} \lambda = \operatorname{argmin}_{\alpha} \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \alpha_n \alpha_n t_n t_n x_n^T x_n \\ \text{under: } \alpha_n \geq 0, \sum_{n=1}^N \alpha_n t_n = 0 \end{cases}$$

(This is known as the dual optimization problem.)

There is no explicit solution for λ but, as the objective function is quadratic, this new problem is relatively easy to solve numerically.

3) Making predictions

$$x_{\text{new}} = \operatorname{sign} \left(\sum_{n=1}^N \alpha_n t_n x_n^T x_{\text{new}} + b \right)$$

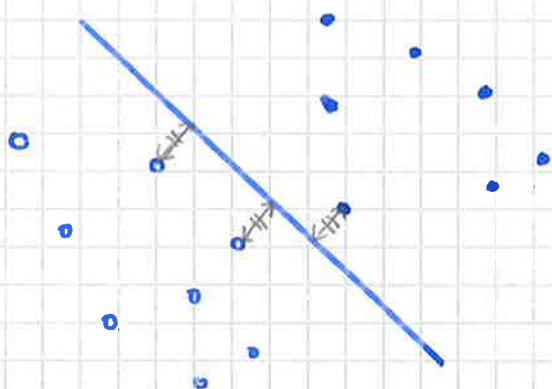
To find b , we use the fact that for the closest point x_n :

$$t_n (w^T x_n + b) = 1$$

and using that $t_n = 1/k_n$, $w = \sum_{n=1}^N \alpha_n t_n x_n$

$$b = t_n - \sum_{n=1}^N \alpha_n t_n x_n^T x_n \quad (b)$$

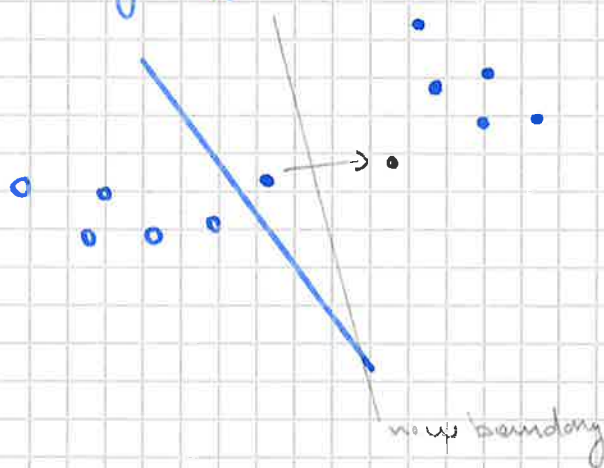
4) Support vectors



The set of points closest to the maximum margin decision boundary are known collectively as the support vectors.

* The decision is function of only a small subset of the training example (we say we have a sparse solution). Consider classifying a test point using KNN when the training data consists of several thousand objects \rightarrow it could take a lot of time

* Drawback: moving a single training point might have a huge effect on the boundary



this method is implementable only if there is a boundary separating the two populations (this is called a hard margin SVM).

5) Soft margins

We want to allow points to potentially lie on the wrong side of the boundary. To achieve this, the constraint becomes

$$k_n(w^T x_n + b) \geq 1 - \xi_n \quad (0 \leq \xi_n \leq 1)$$

New optimization task :

$$\begin{cases} \operatorname{argmin}_w \frac{1}{2} w^T w + C \sum_{n=1}^N \xi_n \\ \text{subject to } \xi_n \geq 0 \text{ and } k_n(w^T x_n + b) \geq 1 - \xi_n \text{ for all } n \end{cases}$$

The new parameter C controls to what extent we are willing to allow points to sit within the margin band or on the wrong side of the decision boundary.

We now need to find the maximum of the following quadratic programming problem :

$$\begin{cases} \operatorname{argmax}_\alpha \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n_1, n_2=1}^N \alpha_{n_1} \alpha_{n_2} k_{n_1} k_{n_2} x_{n_1}^T x_{n_2} \\ \text{subject to } \sum_{n=1}^N \alpha_n k_n = 0 \text{ and } 0 \leq \alpha_n \leq C \text{ for all } n \end{cases}$$

The constant C needs to be fixed. We can set this using cross-validation.

To compute b , we need to find the support vector with the highest value of $w^T x_n$ (or $\sum_n \alpha_n k_n x_n^T x_n$) and compute b from equation (6)

TP: p. 56-60 of the O'Reilly book