

COQUELIN Pierre

CHERON Guilhem

MVA "Kernel methods"

HW 4

Exercise 1

$$a) \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^m \log(1 + e^{-y_i \beta^T x_i}) + \lambda \|\beta\|_2^2 \right\}, \quad x_i \in \mathbb{R}^p, \quad y_i \in \{-1, 1\}$$

$$\begin{aligned} \log(1 + \exp(-y_i \beta^T x_i)) &= \log(1 + \exp(-y_i \log\left(\frac{P_\beta(Y=1|X=x_i)}{P_\beta(Y=-1|X=x_i)}\right))) \\ &= \log\left(1 + \left(\frac{P_\beta(Y=1|X=x_i)}{P_\beta(Y=-1|X=x_i)}\right)^{-y_i}\right) \end{aligned}$$

$$\text{but, } \left(\frac{P_\beta(Y=1|X=x_i)}{P_\beta(Y=-1|X=x_i)}\right)^{-y_i} = \begin{cases} \frac{P_\beta(Y=-1|X=x_i)}{P_\beta(Y=1|X=x_i)} & \text{if } y_i = 1 \\ \frac{P_\beta(Y=1|X=x_i)}{P_\beta(Y=-1|X=x_i)} & \text{if } y_i = -1 \end{cases}$$

$$= \frac{P_\beta(Y=-y_i|X=x_i)}{P_\beta(Y=y_i|X=x_i)}$$

$$= \frac{1 - P_\beta(Y=y_i|X=x_i)}{P_\beta(Y=y_i|X=x_i)}$$

$$= \frac{1 - P_\beta(Y=y_i|X=x_i)}{P_\beta(Y=y_i|X=x_i)}$$

then,

$$\begin{aligned} \log(1 + \exp(-y_i \beta^T x_i)) &= \log\left(1 + \frac{1 - P_\beta(Y=y_i|X=x_i)}{P_\beta(Y=y_i|X=x_i)}\right) \\ &= \log\left(\frac{P_\beta(Y=y_i|X=x_i) + 1 - P_\beta(Y=y_i|X=x_i)}{P_\beta(Y=y_i|X=x_i)}\right) \end{aligned}$$

$$= -\log(P_\beta(Y=y_i|X=x_i))$$

Then we have,

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^m -\log(P_\beta(Y=y_i|X=x_i)) + \lambda \|\beta\|_2^2 \right\}, \text{ which is}$$

$$\text{equivalent to } \max_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^m \log(P_\beta(Y=y_i|X=x_i)) - \lambda \|\beta\|_2^2 \right\}$$

$$b) \quad \beta^T x_i = \sum_{j=1}^m \alpha_j K(x_i, x_j) \quad \Rightarrow \quad \begin{pmatrix} \beta^T x_1 \\ \vdots \\ \beta^T x_m \end{pmatrix} = K \alpha$$

$$\|\beta\|_2^2 = \alpha^T K \alpha \quad (\text{Representer theorem with linear kernel})$$

Then we get the new problem,

$$\min_{\alpha \in \mathbb{R}^m} \{ R(K\alpha) + \lambda \alpha^T K \alpha \}$$

with $R(x) = \sum_{i=1}^m \log(1 + \exp(-y_i x_i))$, here we see that we find the closed and convex risk of the kernel logistic regression.

Let $\alpha^* \in \mathbb{R}^m$ be the solution of this problem.

If we choose K the kernel Gram matrix with entries $K_{ij} = x_i^T x_j$ then we get:

$$\beta^T x_i = \sum_{j=1}^m \alpha_j x_j^T x_i \quad \Rightarrow \quad \beta = \sum_{i=1}^m \alpha_i x_i$$

Then the solution of the first problem is:

$$\beta^* = \sum_{i=1}^m \alpha_i^* x_i$$

c) Rewriting the problem as:

$$\min_{\alpha \in \mathbb{R}^n, u \in \mathbb{R}^n} \{ R(u) + \lambda \alpha^T K \alpha \} \text{ s.t. } K\alpha = u$$

Its Lagrangian is,

$$L(\alpha, u, v) = R(u) + \lambda \alpha^T K \alpha + v^T (K\alpha - u)$$

Thus we have,

$$g(v) = \inf_{\alpha, u} L(\alpha, u, v) = \inf_{\alpha, u} \{ R(u) + \lambda \alpha^T K \alpha + v^T (K\alpha - u) \}$$

We have

$$\inf_u (R(u) - v^T u) = -\sup_u (v^T u - R(u)) = -R^*(v)$$

Taking the gradient w.r.t α and setting it to zero:

$$\nabla_{\alpha} g(v) = 2\lambda \alpha^T K + v^T K = 0$$

$$\Leftrightarrow \alpha = -\frac{v}{2\lambda} + c, \quad c \in \text{Ker}(K) \text{ and we take } c=0$$

$$\begin{aligned} \text{Then } g(v) &= -R^*(v) + \lambda \times -\frac{v^T}{2\lambda} K \times -\frac{v}{2\lambda} + v^T K \times -\frac{v}{2\lambda} \\ &= -R^*(v) - \frac{v^T K v}{4\lambda} \end{aligned}$$

$$R^*(v) = \sup_{x \in \mathbb{R}^n} \{ \langle x, v \rangle - R(x) \}$$

$$= \sup_{x \in \mathbb{R}^n} \left\{ \sum_{i=1}^n x_i v_i - \log(1 + \exp(-y_i x_i)) \right\}$$

$$= \sum_{i=1}^n \sup_{x_i \in \mathbb{R}} \{ x_i v_i - \log(1 + \exp(-y_i x_i)) \}$$

$$= \sum_{i=1}^n \sup_{x_i \in \mathbb{R}} \{ x_i v_i - \ell(y_i, x_i) \}$$

$$= \sum_{i=1}^n \ell_L^*(y_i, v_i), \quad \text{with } \ell_L \text{ which is the logistic}$$

loss and ℓ_L^* is its Fenchel-Legendre transform

We know that,

$$l_L^*(y_i, \mathbf{x}_i) = \begin{cases} (-y_i \mathbf{x}_i) \log(-y_i \mathbf{x}_i) + (1+y_i \mathbf{x}_i) \log(1+y_i \mathbf{x}_i) & \text{if } -1 \leq y_i \mathbf{x}_i \leq 0 \\ +\infty & \text{otherwise} \end{cases}$$

Then we obtain the dual problem

$$\max_{\mathbf{v} \in \mathbb{R}^m} -R^*(\mathbf{v}) - \frac{\mathbf{v}^T \mathbf{K} \mathbf{v}}{4\lambda}$$

$$= \max_{\mathbf{v} \in \mathbb{R}^m} \sum_{i=1}^m (y_i \mathbf{x}_i) \log(-y_i \mathbf{x}_i) - (1+y_i \mathbf{x}_i) \log(1+y_i \mathbf{x}_i) - \frac{\mathbf{v}^T \mathbf{K} \mathbf{v}}{4\lambda}$$

$$\text{s.t. } -1 \leq y_i \mathbf{x}_i \leq 0, i=1, \dots, m$$

If \mathbf{v}^* is the solution of the dual problem, then we get $\alpha^* = -\frac{\mathbf{v}^*}{2\lambda}$

d. * The objective of problem (1) is:

$$C_1(\beta) = \sum_i^n \log(1 + e^{-y_i \beta^T x_i}) + \lambda \|\beta\|_2^2 \\ = -\sum_i^n \log g(y_i \beta^T x_i) + \lambda \|\beta\|_2^2$$

With $g(x) = \frac{1}{1 + \exp(-x)}$ the sigmoid function.

We have

$$\nabla_{\beta} (\log g(y_i \beta^T x_i)) = \frac{1}{g(y_i \beta^T x_i)} g(y_i \beta^T x_i) (1 - g(y_i \beta^T x_i)) (y_i x_i) \\ = y_i x_i (1 - g(y_i \beta^T x_i))$$

Because $\frac{dg(x)}{dx} = g(x)(1 - g(x))$.

Thus the gradient of $C_1(\beta)$ is

$$\nabla_{\beta} C_1(\beta) = -\sum_i^n y_i x_i (1 - g(y_i \beta^T x_i)) + 2\lambda \beta$$

We have

$$\frac{\partial^2 y_i x_i g(y_i \beta^T x_i)}{\partial \beta_p \partial \beta_k} = x_i(p) (g(y_i \beta^T x_i) (1 - g(y_i \beta^T x_i))) x_i(k) y_i^2$$

Thus the Hessian of $C_1(\beta)$ is

$$(\nabla_{\beta}^2 C_1(\beta))_{pk} = \sum_i^n x_i(p) x_i(k) (g(y_i \beta^T x_i) (1 - g(y_i \beta^T x_i))) \\ + 2\lambda \mathbb{1}_{\{p=k\}}$$

* The objective of problem (2) is

$$C_2(\alpha) = \sum_i \log(1 + \exp(-y_i \sum_j \alpha_j K_{ij})) + \lambda \alpha^T K \alpha$$

$$= - \sum_i \log(g(y_i \sum_j \alpha_j K_{ij})) + \lambda \alpha^T K \alpha$$

We have

$$\frac{\partial \log g(y_i \sum_j \alpha_j K_{ij})}{\partial \alpha_k} = (1 - g(y_i \sum_j \alpha_j K_{ij})) y_i K_{ik}$$

Thus the gradient is

$$(\nabla_{\alpha} C_2(\alpha))_k = - \sum_i (1 - g(y_i \sum_j \alpha_j K_{ij})) y_i K_{ik} + 2\lambda (K\alpha)_k$$

And the Hessian is

$$(\nabla_{\alpha}^2 C_2(\alpha))_{kp} = \sum_i g(y_i \sum_j \alpha_j K_{ij}) (1 - g(y_i \sum_j \alpha_j K_{ij})) K_{ik} K_{ip} + 2\lambda K_{kp}$$

* The objective of problem (3) is

$$C_3(r) = \sum_i y_i r_i \log(-y_i r_i) - (1 + y_i r_i) \log(1 + y_i r_i) - \frac{r^T K r}{4\lambda}$$

We have

$$\frac{\partial C_3(r)}{\partial r_i} = y_i \log(-y_i r_i) + y_i - y_i \log(1 + y_i r_i) - y_i - \frac{(K r)_i}{2\lambda}$$

$$= y_i \log\left(\frac{-y_i r_i}{1 + y_i r_i}\right) - \frac{(K r)_i}{2\lambda}$$

and

$$(\nabla_r^2 C_3(r))_{ij} = 1_{\{i=j\}} \left(\frac{1}{y_i r_i} - \frac{1}{1 + y_i r_i} \right) - \frac{K_{ij}}{2\lambda}$$

We implemented a simple Newton-Raphson algorithm and the gradient and Hessian matrix of each problem in MATLAB.

Data is generated randomly using Gaussian distributions with random covariance matrix generated using Wishart distribution. The data is split into a train set and a test set.

The parameter λ is determined using a 5-fold cross validation on the training set. For optimization problems (2) and (3) we used both a linear kernel and a Gaussian kernel with $\sigma = 10$.

Problem (3) is the only one with constraints. After a Newton-Raphson iteration, we set the variables v_i that don't respect these constraints at the closest frontier of the constraint set.

The figures below display the objective function and prediction error in function of Newton-Raphson iterations. Each figure display the results obtained with the various methods, for given values of n and p .

We observe the following points:

- As expected, results obtained with problem (1) and (2) are the same when using a linear kernel (red and yellow curves on the figures below are superimposed).
- Optimization problem (2) needs fewer iterations to converge than problem (3)
- The highest computation cost source is the inversion of the Hessian matrix. Thus, when the number of points n is superior to the dimension p , problem (2) and (3) have a higher computation cost than problem (1). Indeed Hessians of problems (2) and (3) have the same size as the Gram matrix (number of points) while Hessian of problem (1) size's is $p \times p$.

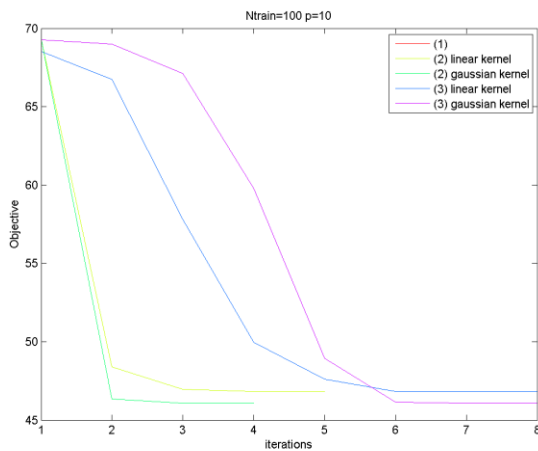


Figure 1 Objective in function of iterations for $n=100$, $p=10$

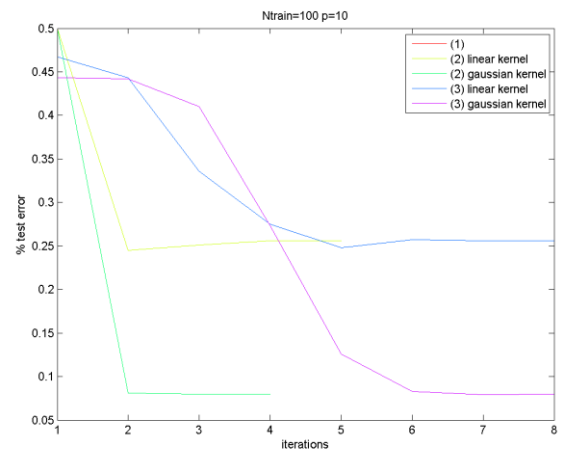


Figure 2 Percentage of error on the test set for $n=100$, $p=10$

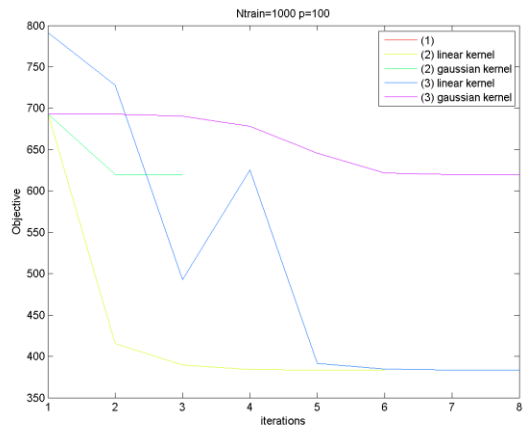


Figure 3 Objective in function of iterations for $n=1000$, $p=100$

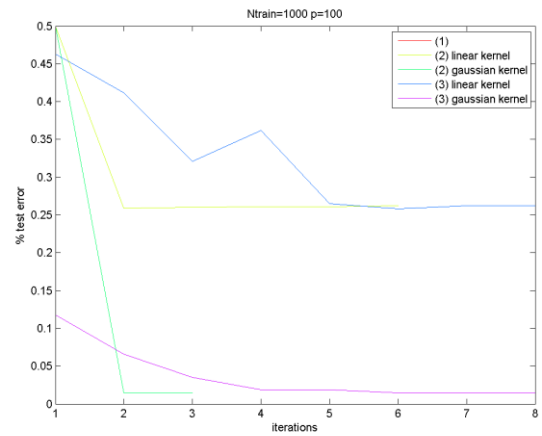


Figure 4 Percentage of error on the test set for $n=1000$, $p=100$