

Problem 1 Generalized Linear Models

(26 points)

We have already introduced linear regression, logistic regression, and multinomial logistic regression. Now we discuss a broader family of models - Generalized Linear Models (GLMs).

1.1 We begin with defining a special family of distributions - the exponential family distributions. If a distribution can be written in the form

Gaussian - $\mathcal{N}(\mu; \sigma)$

$$p(y; \eta) = b(y) \exp(\eta^T t(y) - a(\eta)), \quad (1)$$

it then belongs to the exponential family. In this problem, we always have $y, b(y), a(\eta)$ as scalars. $\eta, t(y) \in \mathbb{R}^K$ are K -dim vectors, but the definition also applies to $K = 1$.

We can easily find that the Bernoulli distribution $y \sim \text{Bernoulli}(q) \Rightarrow p(y) = q^y(1-q)^{1-y}$ is in the exponential family:

$\int b(y) \exp(\eta^T t(y) - a(\eta)) = 1$

$$\begin{aligned} p(y) &= \exp(y \log q + (1-y) \log(1-q)) \\ &= 1 \cdot \exp(\log \frac{q}{1-q} \cdot y + \log(1-q)), \text{ where} \\ b(y) &= 1 \\ \eta &= \log \frac{q}{1-q} \\ t(y) &= y \\ a(\eta) &= -\log(1-q). \end{aligned}$$

$y \log q + \log(1-q) - y \log(1-q)$

Show that the categorical distribution is also in the exponential family, and write down its $b(y), \eta, t(y), a(\eta)$. (8 points)

(Hint: You may consider using the following form of categorical distribution:

$$p(y; q) = (Cq_1)^{1\{y=1\}} (Cq_2)^{1\{y=2\}} \dots (Cq_K)^{1\{y=K\}},$$

where q_k is a non-negative scalar. $C = 1 / \sum_{k=1}^K q_k$ for normalization so that Cq_1, \dots, Cq_K are probabilities. $1\{y=k\} = 1$ if $y=k$, otherwise $1\{y=k\} = 0$. Notice that η can be some expression of q .

$q_1 \sim q_K \quad K-1$

$$\begin{aligned} p(y; q) &= \exp\left(\sum_{k=1}^K 1\{y=k\} \log Cq_k\right) \\ &= \exp\left(\sum_{k=1}^K 1\{y=k\} \log q_k + \sum_{k=1}^K 1\{y=k\} \log C\right) \\ &= \exp(\eta^T e_y + \log C), \end{aligned}$$

$1 \cdot \log C$
 $[1\{y=1\} \dots 1\{y=K\}] = t(y)$
 $[\log q_1 \dots \log q_K] = \eta$

where $\eta = (\log q_1, \log q_2, \dots, \log q_K)$, and e_y is a K -dim one-hot vector (only the y -th element is 1 and

the others are 0). Therefore, the categorical distribution is in the exponential family, and

$$\begin{aligned} b(y) &= 1, \\ \boldsymbol{\eta} &= (\log q_1, \log q_2, \dots, \log q_K), \\ \mathbf{t}(y) &= \mathbf{e}_y, \\ a(\boldsymbol{\eta}) &= \log\left(\sum_{k=1}^K q_k\right). \end{aligned}$$

(4 points points)

1.2 Now we give the steps to construct a GLM:

- (1) Given the input feature $\mathbf{x} \in \mathbb{R}^D$, find a proper distribution belonging to the exponential family as the distribution of the label y conditioning on \mathbf{x} : $p(y|\mathbf{x}) \sim \text{ExponentialFamily}(\boldsymbol{\eta})$. $\boldsymbol{\eta} \in \mathbb{R}^K$.
- (2) To make the model linear, we let $\boldsymbol{\eta} = \mathbf{W}\mathbf{x}$. $\mathbf{W} \in \mathbb{R}^{K \times D}$. (When $K = 1$ we usually write it as $w^T \mathbf{x}$.)
- (3) We select $h(\mathbf{x}; \mathbf{W}) = \mathbb{E}_{y \sim p(y|\mathbf{x}; \mathbf{W})} \mathbf{t}(y)$ as our predicted value.

If we select the conditional distribution in Step (1) as the Bernoulli distribution, please finish the remaining steps to construct a GLM and show $h(\mathbf{x}; \mathbf{w})$. (6 points)

$$\eta = \log \frac{q}{1-q} \Rightarrow q = \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\eta}}.$$

With $\eta = w^T \mathbf{x}$, we predict

$$h(\mathbf{x}; \mathbf{w}) = \mathbb{E}(t(y)) = q \cdot 1 + (1-q) \cdot 0 = q = \frac{1}{1+e^{-w^T \mathbf{x}}}.$$

We find that the Bernoulli distribution assumption on the label leads to logistic regression.

1.3 If we select the conditional distribution in Step (1) as the categorical distribution in the previous question, please finish the remaining steps to construct a GLM and show $h(\mathbf{x}; \mathbf{W})$. (6 points)

Let $\boldsymbol{\eta} = \mathbf{W}\mathbf{x}$, we have $\log q_k = w_k^T \mathbf{x} \Rightarrow q_k = \exp(w_k^T \mathbf{x})$, $\forall k = 1, 2, \dots, K$. $C = 1 / \sum_{k=1}^K \exp(w_k^T \mathbf{x})$.

$$\begin{aligned} h(\mathbf{x}; \mathbf{W}) &= \mathbb{E}(\mathbf{t}(y)) = \sum_{k=1}^K P(y = k; \mathbf{q}) \mathbf{e}_y \\ &= \sum_{k=1}^K C q_k \mathbf{e}_y \quad \left(\frac{q_k}{\sum_i q_i} \right) \\ &= \left(\frac{\exp(w_1^T \mathbf{x})}{\sum_{k=1}^K \exp(w_k^T \mathbf{x})}, \dots, \frac{\exp(w_K^T \mathbf{x})}{\sum_{k=1}^K \exp(w_k^T \mathbf{x})} \right) \leftarrow \\ &= \text{Softmax}(\mathbf{W}\mathbf{x}). \end{aligned}$$

Handwritten notes:
 $\log q_k = \eta_k = w_k^T \mathbf{x}$
 $\eta = \mathbf{W}\mathbf{x}$

We can find that by taking the assumption that the label follows a categorical distribution, the derived GLM is exactly the (randomized) multinomial logistic regression.

1.4 Now let's construct a GLM to predict values that are most likely to follow the Poisson distribution (e.g. daily number of visitors in a store). Show that Poisson distribution is in the exponential family and finish the steps to construct a GLM by deriving $h(x; w)$. (A slight difference in these steps is that η in this question is now a scalar.) **(6 points)**

(Hint: You may consider using the following form of Poisson distribution:

$$p(y; \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!},$$

where $\lambda > 0$ and is a scalar. Notice that η can be some expression of λ .)

$$p(y; \lambda) = \frac{1}{y!} \exp(y \log \lambda - \lambda).$$

Therefore,

$$b(y) = 1/y!$$

$$\eta = \log \lambda = w^T x$$

$$t(y) = y$$

$$a(\eta) = \lambda.$$

Let $\eta = w^T x$, we have $\lambda = \exp(\eta) = \exp(w^T x)$. Therefore,

$$h(x; w) = \mathbb{E}(t(y)) = \exp(-\lambda) \sum_{y=0}^{+\infty} \frac{y \lambda^y}{y!}$$

$$= \exp(-\lambda) \sum_{y=1}^{+\infty} \frac{\lambda^{y-1}}{(y-1)!} \cdot \lambda$$

$$= \exp(-\lambda) \sum_{y=0}^{+\infty} \frac{\lambda^y}{(y)!} \cdot \lambda$$

$$= \exp(-\lambda) \exp(\lambda) \exp(w^T x) = \exp(w^T x).$$

Problem 2 Neural Networks

In the lecture, we have talked about error-backpropagation, a way to compute partial derivatives (or gradients) w.r.t the parameters of a neural network to optimize using gradient descent. In this question, you are going to practice (Q2.1) error-backpropagation, (Q2.2) how initialization affects optimization, and (Q2.3) the importance of nonlinearity.

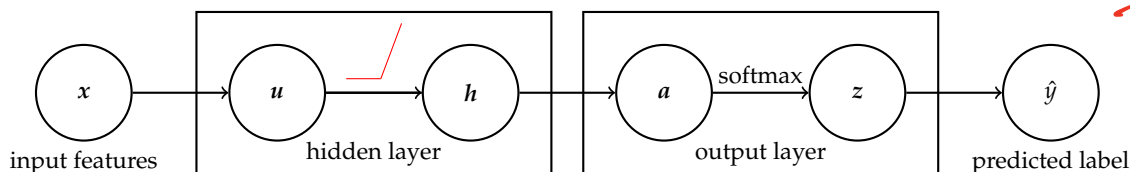


Figure 1: A diagram of a 1-hidden layer neural net. The edges mean mathematical operations, and the circles mean variables. Generally we call the combination of a linear (or affine) operation and a nonlinear operation (like element-wise sigmoid or the rectified linear unit (relu) operation as in eq. (4)) as a hidden layer. Note the two slight differences compared to the diagram used in the lecture : 1) one circle represents a vector and thus an array of neurons here and 2) the activation operations are also explicitly represented as edges here.

Specifically, you are given the following 1-hidden layer neural net for a K -class classification problem (see Fig. 1 for illustration and details), and $(\mathbf{x} \in \mathbb{R}^D, y \in \{1, 2, \dots, K\})$ is a labeled instance,

$$\mathbf{x} \in \mathbb{R}^D \quad (2)$$

$$\mathbf{u} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}, \quad \mathbf{W}^{(1)} \in \mathbb{R}^{M \times D} \text{ and } \mathbf{b}^{(1)} \in \mathbb{R}^M \quad (3)$$

$$\mathbf{h} = \max\{0, \mathbf{u}\} = \begin{bmatrix} \max\{0, u_1\} \\ \vdots \\ \max\{0, u_M\} \end{bmatrix} \quad (4)$$

$$\mathbf{a} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}, \quad \mathbf{W}^{(2)} \in \mathbb{R}^{K \times M} \text{ and } \mathbf{b}^{(2)} \in \mathbb{R}^K \quad (5)$$

$$\mathbf{z} = \begin{bmatrix} e^{a_1} \\ \frac{e^{a_1}}{\sum_k e^{a_k}} \\ \vdots \\ e^{a_K} \\ \frac{e^{a_K}}{\sum_k e^{a_k}} \end{bmatrix} \quad (6)$$

$$\hat{y} = \arg \max_k z_k. \quad (7)$$

For K -class classification problem, one popular loss function for training is the cross-entropy loss. Specifically we denote the cross-entropy loss with respect to the training example (\mathbf{x}, y) by l :

$$l = -\ln(z_y) = \ln \left(1 + \sum_{k \neq y} e^{a_k - a_y} \right) = \sum_i e^{a_i - a_y} = \frac{\sum e^{a_i}}{e^{a_y}}$$

Note that l is a function of the parameters of the network, that is, $\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}$ and $\mathbf{b}^{(2)}$.

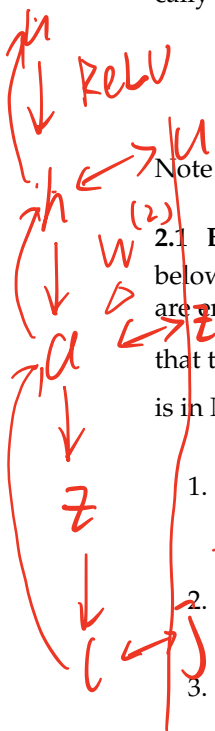
2.1 Error Back-propagation Assume that you have computed $\mathbf{u}, \mathbf{h}, \mathbf{a}, \mathbf{z}$, given (\mathbf{x}, y) . Follow the four steps below to find out the derivatives of l with respect to all the four parameters $\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}$ and $\mathbf{b}^{(2)}$. You are encouraged to use matrix/vector forms to simplify your answers. Note that we follow the convention that the derivative with respect to a variable is of the same dimension of that variable. For example, $\frac{\partial l}{\partial \mathbf{W}^{(1)}}$ is in $\mathbb{R}^{M \times D}$. (This is called the **denominator layout**.)

1. First express $\frac{\partial l}{\partial \mathbf{a}}$ in terms of \mathbf{z} and y . You may find it convenient to use the notation $\mathbf{y} \in \mathbb{R}^K$ whose k -th coordinate is 1 if $k = y$ and 0 otherwise. (4 points)

2. Then express $\frac{\partial l}{\partial \mathbf{W}^{(2)}}$ and $\frac{\partial l}{\partial \mathbf{b}^{(2)}}$ in terms of $\frac{\partial l}{\partial \mathbf{a}}$ and \mathbf{h} . (4 points)

3. Next express $\frac{\partial l}{\partial \mathbf{u}}$ in terms of $\frac{\partial l}{\partial \mathbf{a}}, \mathbf{u}$, and $\mathbf{W}^{(2)}$. You will need to use the (sub)derivative of the ReLU function $\max\{0, u\}$ denoted by $H(u)$ and is 1 if $u \geq 0$ and 0 otherwise. Also, you may find it convenient to use the notation $\mathbf{H}(u) \in \mathbb{R}^{M \times M}$ which stands for a diagonal matrix with $H(u_1), \dots, H(u_M)$ on the diagonal. (4 points)

4. Finally, express $\frac{\partial l}{\partial \mathbf{W}^{(1)}}$ and $\frac{\partial l}{\partial \mathbf{b}^{(1)}}$ in terms of $\frac{\partial l}{\partial \mathbf{u}}$ and \mathbf{x} . (4 points)



$$\frac{\partial l}{\partial \mathbf{u}}$$

$$\frac{\partial l}{\partial \mathbf{h}} = \mathbf{W}^{(2)T} \frac{\partial l}{\partial \mathbf{a}}$$

$$\frac{\partial l}{\partial \mathbf{u}} = \mathbf{W}^{(2)T} \frac{\partial l}{\partial \mathbf{a}} \odot \mathbf{H}'(\mathbf{u})$$

$$\frac{\partial l}{\partial a_i} = \frac{1}{1 + \sum_{k \neq y} e^{a_k - a_y}} \cdot (a_i - a_y) \cdot (-1)$$

$$\frac{\partial l}{\partial \mathbf{a}} = \mathbf{z} - \mathbf{y}$$

$$\frac{\partial l}{\partial \mathbf{W}^{(2)}} = \mathbf{h} (\mathbf{z} - \mathbf{y})$$

$$\frac{\partial l}{\partial \mathbf{b}^{(2)}} = \mathbf{z} - \mathbf{y}$$

$\frac{\partial L}{\partial a}, W^{(2)}, H(u) \Rightarrow i)$

$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \otimes \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \Leftrightarrow \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} x \\ \cdot \\ \cdot \end{bmatrix}$

$\frac{\partial L}{\partial u_i}$

$\Delta \frac{\partial L}{\partial a} = z - y$

$\Delta \frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial a} h^T$

$\Delta \frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial a}$

$\Delta \frac{\partial L}{\partial u} = \frac{\partial h}{\partial u} \frac{\partial L}{\partial h} \frac{\partial L}{\partial a} = H(u) W^{(2)T} \frac{\partial L}{\partial a}$

$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial u} x^T$

$\frac{\partial L}{\partial b^{(1)}} = \frac{\partial L}{\partial u}$

$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial u} \Rightarrow m \times m$

$\frac{\partial h}{\partial u} \quad m \times k \quad \frac{\partial L}{\partial a} \quad k \times 1$

$(W^{(2)})^T \quad m \times 1$

(derivation is in Slide 17 of Lec 4)

2.2 Initialization Suppose we initialize $W^{(1)}, W^{(2)}, b^{(1)}$ with zero matrices/vectors (i.e., matrices and vectors with all elements set to 0), please first verify that $\frac{\partial L}{\partial W^{(1)}}, \frac{\partial L}{\partial W^{(2)}}, \frac{\partial L}{\partial b^{(1)}}$ are all zero matrices/vectors, irrespective of x, y and the initialization of $b^{(2)}$.

Now if we perform stochastic gradient descent for learning the neural network, please explain with a concise statement why no learning will happen with this initialization. **(4 points)**

Since $W^{(2)}$ is all zero, $\frac{\partial L}{\partial a}$ is all zero. So $\frac{\partial L}{\partial W^{(1)}}, \frac{\partial L}{\partial b^{(1)}}$ are all zero. Since $W^{(1)}, b^{(1)}$ are all zero, h is all zero. So $\frac{\partial L}{\partial W^{(2)}}$ is all zero. In each iteration, all gradients with respect to these three parameters are zero, so no updates will be made.

2.3 Non-linearity As mentioned in the lecture, non-linearity is very important for neural networks. With non-linearity (e.g., eq. (4)), the neural network shown in Fig. 1 can be seen as a nonlinear basis function ϕ (i.e., $\phi(x) = h$) followed by a linear classifier f (i.e., $f(h) = y$).

Please show that, by removing the nonlinear operation in eq. (4) and setting eq. (5) to be $a = W^{(2)}u + b^{(2)}$, the resulting network is essentially a linear classifier. More specifically, you can now represent a as $Ux + v$, where $U \in \mathbb{R}^{K \times D}$ and $v \in \mathbb{R}^K$. Please write down the representation of U and v using $W^{(1)}, W^{(2)}, b^{(1)}$, and $b^{(2)}$. **(4 points)**

By combining the equations, we can get:

$a = W^{(2)}u + b^{(2)}$

$= W^{(2)}(W^{(1)}x + b^{(1)}) + b^{(2)}$

$= (W^{(2)}W^{(1)})x + (W^{(2)}b^{(1)} + b^{(2)})$

$U = W^{(2)}W^{(1)}$

$v = W^{(2)}b^{(1)} + b^{(2)}$

Problem 3 Regularized Linear Regression With Kernels (15 points)

In class, we derive the closed-form solution of regularized linear regression with kernels. Now we discuss its gradient descent solution.

For the following regularized linear regression with feature mapping $\phi \in \mathbb{R}^D \rightarrow \mathbb{R}^M, M \gg D$

$$L(w) = \frac{1}{2} \sum_{i=1}^n \|w^T \phi(x_i) - y_i\|_2^2 + \frac{1}{2} \lambda \|w\|_2^2, \lambda > 0,$$

3.1 Write down w_{t+1} after one step gradient descent (using all examples) from w_t with learning rate $\alpha > 0$. (3 points)

$$\nabla_{w_t} L = \sum_{i=1}^n (w_t^T \phi(x_i) - y_i) \phi(x_i) + \lambda w_t$$

(1 points)

$$w_{t+1} = w_t - \alpha \left(\sum_{i=1}^n (w_t^T \phi(x_i) - y_i) \phi(x_i) + \lambda w_t \right)$$

$$= (1 - \alpha \lambda) w_t + \alpha \left(\sum_{i=1}^n (y_i - w_t^T \phi(x_i)) \phi(x_i) \right)$$

(2 points)

3.2 What will be the problem if we directly conduct gradient descent? (2 points)

Since $M \gg D$, directly calculating the gradient has high computation and memory costs. It is even not applicable when M is infinity.

3.3 Denote as K the corresponding kernel of ϕ : $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$.

(1) Prove that if we start from $w_0 = 0$, for each t during gradient descent, we can always find scalars $\beta_i^{(t)}$, $i = 1, \dots, n$ such that $w_t = \sum_{i=1}^n \beta_i^{(t)} \phi(x_i)$. In other words, each w_t is a linear combination of $\phi(x_1), \dots, \phi(x_n)$. (Hint: use induction from $t = 0$ to $1, 2, \dots$) (8 points)

(2) Write down $\beta_1^{(t+1)}, \dots, \beta_n^{(t+1)}$ after one step gradient descent from $\beta_1^{(t)}, \dots, \beta_n^{(t)}$. Note that you should not have w in your final result. (2 points)

(1) • Trivial case: $t = 0$. We can simply find $\beta_i^{(0)} = 0$, $i = 1, \dots, n$.

• Induction: If $w_t = \sum_{i=1}^n \beta_i^{(t)} \phi(x_i)$,

$$\begin{aligned} w_{t+1} &= (1 - \alpha \lambda) w_t + \alpha \left(\sum_{i=1}^n (y_i - w_t^T \phi(x_i)) \phi(x_i) \right) \\ &= \sum_{i=1}^n \left((1 - \alpha \lambda) \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} \phi(x_j)^T \phi(x_i)) \right) \phi(x_i) \\ &= \sum_{i=1}^n \left((1 - \alpha \lambda) \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} K(x_i, x_j)) \right) \phi(x_i) \\ &\quad \underbrace{\left((1 - \alpha \lambda) \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} K(x_i, x_j)) \right)}_{\beta_i^{(t+1)}} \phi(x_i) \end{aligned}$$

(2)

$$\beta_i^{(t+1)} = (1 - \alpha \lambda) \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} K(x_i, x_j)), \quad i = 1, 2, \dots, n$$

Problem 4 Direction of Linear Discriminant Hyperplane

(15 points)

Consider linear discriminant analysis for a two-class classification problem on a dataset of N inputs $\{\mathbf{x}_1 \dots \mathbf{x}_N\}$ and corresponding labels $\{y_1 \dots y_N\}$, $y_i \in \{-1, 1\} \forall i \in \{1 \dots N\}$. We say input \mathbf{x}_i belongs to class \mathcal{C}_1 if its label y_i is 1 and it belongs to class \mathcal{C}_{-1} if its label is -1. Mathematically, $\mathcal{C}_1 = \{(\mathbf{x}_i, y_i) : i \in [N], y_i = 1\}$ and $\mathcal{C}_{-1} = \{(\mathbf{x}_i, y_i) : i \in [N], y_i = -1\}$

We aim to find a separating hyperplane \mathbf{w} such that if input \mathbf{x}_i belongs to \mathcal{C}_1 then $\mathbf{w}^T \mathbf{x}_i \geq 0$ and if it belongs to \mathcal{C}_{-1} then $\mathbf{w}^T \mathbf{x}_i \leq 0$. However, this might not be always possible. Instead, one way to relax the goal is to find a hyperplane \mathbf{w}^* that maximizes $f(\mathbf{w}) = \sum_{i=1}^N y_i \mathbf{w}^T \mathbf{x}_i$ under the constraint $\|\mathbf{w}\| = 1$. Note that $f(\mathbf{w})$ can be arbitrarily maximized by increasing the magnitude of \mathbf{w} and thus the constraint $\|\mathbf{w}\| = 1$ (or equivalently, $\|\mathbf{w}\|^2 = 1$) is important. We also assume that $\sum_{i=1}^N y_i \mathbf{x}_i \neq \mathbf{0}$ otherwise the objective $f(\mathbf{w})$ is always 0.

This can be written as a well-defined optimization problem using Lagrange multipliers (you do not have to know what this is to solve this problem). More concretely, there exists $\lambda \neq 0$ such that the hyperplane \mathbf{w}^* we are looking for satisfies:

$$\mathbf{w}^* = \arg \max_{\mathbf{w} \in \mathbb{R}^D} \sum_{i=1}^N y_i \mathbf{w}^T \mathbf{x}_i - \lambda (\mathbf{w}^T \mathbf{w} - 1) \quad (8)$$

4.1 Prove the following

(8 points)

primal: $\arg \min_{\mathbf{w} \in \mathbb{R}^D} \sum_{i=1}^N -y_i \mathbf{w}^T \mathbf{x}_i$, s.t. $\|\mathbf{w}\| = 1$

dual: $\arg \max_{\lambda} \sum_{i=1}^N -y_i \mathbf{w}^T \mathbf{x}_i + \lambda (\mathbf{w}^T \mathbf{w} - 1)$

To find the maximum we set the gradient of $f(\mathbf{w}) = \sum_{i=1}^N y_i \mathbf{w}^T \mathbf{x}_i + \lambda (\mathbf{w}^T \mathbf{w} - 1)$ to 0.

$$\nabla f(\mathbf{w}) = \sum_{i=1}^N y_i \mathbf{x}_i - 2\lambda \mathbf{w} = \mathbf{0}$$

$$\Rightarrow \mathbf{w}^* = \frac{1}{2\lambda} \left(\sum_{i=1}^N y_i \mathbf{x}_i \right) = \frac{1}{2\lambda} \left(\sum_{i: \mathbf{x}_i \in \mathcal{C}_1} \mathbf{x}_i - \sum_{j: \mathbf{x}_j \in \mathcal{C}_{-1}} \mathbf{x}_j \right)$$

4.2 Find the value of λ .

(4 points)

Since $\|\mathbf{w}^*\| = 1$ we know $\lambda = \frac{1}{2} \left\| \sum_{i: \mathbf{x}_i \in \mathcal{C}_1} \mathbf{x}_i - \sum_{j: \mathbf{x}_j \in \mathcal{C}_{-1}} \mathbf{x}_j \right\|$.

4.3 In terms of minimizing the training error, can you think of one issue of our objective, i.e. maximizing $f(\mathbf{w})$? (3 points)

Maximizing this objective might lead to a solution that prefers having a large margin on some data points with the price of misclassifying others.

$$\text{primal: } \star \arg\min_w \sum_i -y_i w^T x_i$$

$$\text{s.t. } w^T w = 1$$

$$\downarrow$$

$$\text{dual: } \arg\max_{\lambda} \arg\min_{w \in \mathbb{R}^p} \sum_i -y_i w^T x_i + \lambda (w^T w - 1)$$

$$\lambda \in \mathbb{R}$$

$$\sum_i -y_i x_i + 2\lambda w^* = 0 \Rightarrow \lambda \neq 0$$

$$w^* = \frac{1}{2\lambda} \sum_i y_i x_i$$

$$\star L(\lambda) = \sum_i -y_i \left(\frac{1}{2\lambda} \sum_j y_j x_j^T \right) x_i$$

$$L(\lambda) = + \lambda \cdot \left[\frac{1}{4\lambda^2} \left(\sum_i y_i x_i^T \right) \left(\sum_j y_j x_j \right) - 1 \right]$$

$$= \frac{1}{2\lambda} \sum_{i,j} -y_i y_j x_j^T x_i$$

$$+ \frac{1}{4\lambda} \sum_{i,j} y_i y_j x_i^T x_j - \lambda$$

$$= - \left(\frac{1}{4\lambda} \sum_{i,j} y_i y_j x_i^T x_j \right) + \lambda$$

$$\sum_i -y_i w^T x_i$$

$$w^* = \frac{1}{2\lambda} \sum_i y_i x_i$$

$$\sum_i \sum_j \frac{1}{2\lambda} y_i y_j x_i^T x_j$$

$$= \frac{1}{2\lambda} \left\| \sum_i y_i x_i \right\|_2^2$$

$$\left(\lambda^* \right)^2 = \frac{1}{4} \left\| \sum_i y_i x_i \right\|_2^2$$

$$\lambda^* = \pm \frac{1}{2} \left\| \sum_i y_i x_i \right\|_2$$

$$x^* = \frac{1}{\left\| \sum_i y_i x_i \right\|_2}$$

ハ・ハ