# Project 3                                    Pooya Taherkhani

## CS 6830 – Machine Learning
Professor Razvan Bunescu

Textbook: PRML by Bishop                        *June 2018*

*Theory &*
*Implementation*

KERNEL TECHNIQUES
KERNEL NEAREST NEIGHBOR
DISTANCE-WEIGHTED NEAREST NEIGHBOR
MATRIX PROPERTIES
MAX MARGIN HYPERPLANES
SVM REGRESSION
LARGE MARGIN PERCEPTRON
FEATURE SELECTION (DOCUMENT CLASSIFIER)

## 1. Kernel Techniques

*Show that if* $k_1(\mathbf{x}, \mathbf{y})$ *and* $k_2(\mathbf{x}, \mathbf{y})$ *are valid kernel functions, then* $k(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y})k_2(\mathbf{x}, \mathbf{y})$ *is also a valid kernel.*

Assume

$$k_1(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^{\mathrm{T}}\phi(\mathbf{y}) \ ; \quad \phi(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) & \dots & \phi_n(\mathbf{x}) \end{bmatrix}^{\mathrm{T}}$$

$$k_2(\mathbf{x}, \mathbf{y}) = \psi(\mathbf{x})^{\mathrm{T}}\psi(\mathbf{y}) \ ; \quad \psi(\mathbf{x}) = \begin{bmatrix} \psi_1(\mathbf{x}) & \dots & \psi_m(\mathbf{x}) \end{bmatrix}^{\mathrm{T}}$$

We need to show that there exists a vector function $\omega$ such that

$$k(\mathbf{x}, \mathbf{y}) = \omega(\mathbf{x})^{\mathrm{T}}\omega(\mathbf{y})$$

Thus, we want to separate $\mathbf{x}$ and $\mathbf{y}$ in the function $k$ to produce the vector function $\omega$,

$$k(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y})k_2(\mathbf{x}, \mathbf{y})$$
$$= \phi(\mathbf{x})^{\mathrm{T}}\phi(\mathbf{y})\psi(\mathbf{x})^{\mathrm{T}}\psi(\mathbf{y})$$
$$= \left(\phi_1(\mathbf{x})\phi_1(\mathbf{y}) + \dots + \phi_n(\mathbf{x})\phi_n(\mathbf{y})\right)\left(\psi_1(\mathbf{x})\psi_1(\mathbf{y}) + \dots + \psi_m(\mathbf{x})\psi_m(\mathbf{y})\right)$$

We can proceed with our argument regardless of the relative size of $n$ and $m$. That is, it can be the case that $n < m$, $n = m$, or $n > m$. By expanding the last equation above, we have

$$k(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{n}\sum_{j=1}^{m} \phi_i(\mathbf{x})\phi_i(\mathbf{y})\psi_j(\mathbf{x})\psi_j(\mathbf{y})$$

$$= \sum_{i=1}^{n}\sum_{j=1}^{m} \phi_i(\mathbf{x})\psi_j(\mathbf{x})\phi_i(\mathbf{y})\psi_j(\mathbf{y}) \quad \text{scalar multiplication is commutative}$$

Therefore, $\omega$ is

$$\omega(\mathbf{x}) = \begin{bmatrix} \phi_i(\mathbf{x})\psi_j(\mathbf{x}) \end{bmatrix}_{nm \times 1} \qquad \begin{array}{c} \text{for all combinations of } i \text{ and } j \\ \text{where } i = 1, \ \dots, n \text{ and} \\ j = 1, \ \dots, m \end{array}$$

## 2. Kernel Nearest Neighbor [Exercise 6.3, page 320]

Nearest Neighbor Classification:

- find the nearest data point $\mathbf{x}_{n'}$ to the new input $\mathbf{x}$
- then $y(\mathbf{x}) = t_{n'}$

where $y(\mathbf{x})$ is the predicted label for $\mathbf{x}$ and $t_{n'}$ is the true label of the nearest point $\mathbf{x}_{n'}$.

$$\mathbf{x}^{\mathrm{T}} = \begin{bmatrix} x_1 & \dots & x_k \end{bmatrix}_{1 \times k} \qquad \text{new input with unknown label}$$

$$\mathbf{x}_n^{\mathrm{T}} = \begin{bmatrix} x_{n1} & \dots & x_{nk} \end{bmatrix}_{1 \times k} \qquad n^{\mathrm{th}} \text{ example with known label, } n \in 1, \dots, N$$

The distance $d$ (used to find the nearest point), at its simplest form, is defined as the Euclidean distance

$$
\begin{aligned}
d^2 = \|\mathbf{x} - \mathbf{x}_n\|^2 &= (x_1 - x_{n1})^2 + \cdots + (x_k - x_{nk})^2 \\
&= x_1^2 - 2x_1 x_{n1} + x_{n1}^2 + \cdots + x_k^2 - 2x_k x_{nk} + x_{nk}^2 \\
&= \left(x_1^2 + \cdots + x_k^2\right) + \left(x_{n1}^2 + \cdots + x_{nk}^2\right) - 2\left(x_1 x_{n1} + \cdots + x_k x_{nk}\right) \\
&= \mathbf{x}^\mathsf{T}\mathbf{x} + \mathbf{x}_n^\mathsf{T}\mathbf{x}_n - 2\mathbf{x}^\mathsf{T}\mathbf{x}_n
\end{aligned}
$$

each term of which is now a dot product that can be represented by a linear kernel.

$$
\|\mathbf{x} - \mathbf{x}_n\|^2 = k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}, \mathbf{x}_n) \tag{1}
$$

Thus the predicted label, $y(\mathbf{x}) = t_n$, can be computed by finding $n \in \{1, \ldots, N\}$ that returns the smallest distance $\|\mathbf{x} - \mathbf{x}_n\|$, where $N$ is the number of examples with known label.

Now the kernel function $k$ in equation (1) can be replaced with any (nonlinear) valid kernel.

## 3. Distance-weighted Nearest Neighbor – Multi Class

Kernel-based distance-weighted nearest neighbor for $K$ classes, where $K \geqslant 2$.

- Assume $K$ classes for label $t$, that is, $t \in \{C_1, ..., C_K\}$.

- The predicted label for $\mathbf{x}$ is $y(\mathbf{x}) = \underset{t}{\operatorname{argmax}} \sum_{n=1}^{N} \kappa(\mathbf{x}, \mathbf{x}_n)\delta_t(t_n)$

  where $\delta_t(x)$ is Kronecker delta function $\delta_t(x) = \begin{cases} 1 & \text{if } x = t \\ 0 & \text{otherwise} \end{cases}$

The algorithm above works as follows. For each class $C_k$ (where $k = 1, \ldots, K$), construct a vector of 0s and 1s with the length $N$ (number of data points) by assigning 1 to each point $\mathbf{x}_n$ that belongs to class $C_k$ (that is, $t_n = C_k$), and 0 otherwise. we will have $K$ vectors of 1s and 0s, each of length $N$. Also for each point $\mathbf{x}_n$ compute a weight $\kappa(\mathbf{x}, \mathbf{x}_n)$ (or a similarity measure between $\mathbf{x}$ and $\mathbf{x}_n$) and multiply it by the corresponding element (which is 0 or 1) in all vectors. Out of the resulting $K$ vectors, pick the one with the largest sum of elements. The class associated with the selected vector is the predicted class for $\mathbf{x}$.

## 4. Matrix Properties [Exercise 6.24, page 322]

*A diagonal matrix $W$ whose elements satisfy $0 < W_{ii} < 1$ is positive definite.*

A symmetric $n \times n$ real matrix $W$ is said to be positive definite if the quadratic form $z^{\mathsf{T}}Wz$ (which yields a scalar) is strictly positive for every non-zero column vector $z$ of $n$ real numbers.

$$z^{\mathsf{T}}Wz > 0 \quad \forall z \neq 0$$

We have

$$z^{\mathsf{T}}Wz =$$

$$\begin{bmatrix} z_1 & \cdots & z_n \end{bmatrix} \begin{bmatrix} W_{11} & & 0 \\ & \ddots & \\ 0 & & W_{nn} \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} =$$

$$\begin{bmatrix} z_1 W_{11} & \cdots & z_n W_{nn} \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} =$$

$$z_1^2 W_{11} + \cdots + z_n^2 W_{nn} \tag{1}$$

We show that the diagonal form $z_1^2 W_{11} + \cdots + z_n^2 W_{nn}$ is strictly positive. Note that $W_{ii} > 0$ and $z_i^2 \geqslant 0$, which implies $z_i^2 W_{ii} \geqslant 0$ for all $i \in \{1, \ldots, n\}$. However, there exists at least one $k \in \{1, \ldots, n\}$ for which $z_k > 0$ because $z$ is a non-zero vector. Therefore, $z_k^2 W_{kk} > 0$. Thus, expression (1) is strictly greater than 0. $\qquad \square$

*The sum of two positive definite matrices is itself positive definite.*

$$z^{\mathsf{T}}Az > 0 \ , \ z^{\mathsf{T}}Bz > 0 \quad \forall z_{n \times 1} \neq 0$$

where $A$ and $B$ are $n \times n$ matrices (They have to be of the same size because they are to be added together). We can rewrite $z^{\mathsf{T}}Az > 0$ as

$$\begin{bmatrix} z_1 & \cdots & z_n \end{bmatrix} \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} > 0$$

$$\begin{bmatrix} z_1 a_{11} + \cdots + z_n a_{n1} & \cdots & z_1 a_{1n} + \cdots + z_n a_{nn} \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} > 0$$

$$(z_1^2 a_{11} + \cdots + z_1 z_n a_{n1}) + \cdots + (z_1 z_n a_{1n} + \cdots + z_n^2 a_{nn}) > 0$$

The last two lines can be written as

$$\left[\sum_{i=1}^{n} z_i a_{i1} \quad \cdots \quad \sum_{i=1}^{n} z_i a_{in}\right] \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} > 0$$

$$z_1 \sum_{i=1}^{n} z_i a_{i1} + \cdots + z_n \sum_{i=1}^{n} z_i a_{in} > 0$$

$$\sum_{j=1}^{n} \sum_{i=1}^{n} z_j z_i a_{ij} > 0$$

$$\sum_{i=1}^{n} \sum_{j=1}^{n} z_i z_j a_{ij} > 0 \qquad (2)$$

Similarly, we can rewrite the quadratic form $z^{\mathsf{T}} B z > 0$ as

$$\sum_{i=1}^{n} \sum_{j=1}^{n} z_i z_j b_{ij} > 0 \qquad (3)$$

We show that $z^{\mathsf{T}}(A + B)z$ is strictly positive for any non-zero vector $z$ of size $n$. We have

$$z^{\mathsf{T}}(A + B)z =$$

$$\begin{bmatrix} z_1 & \cdots & z_n \end{bmatrix} \left( \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} + \begin{bmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \cdots & b_{nn} \end{bmatrix} \right) \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} =$$

$$\begin{bmatrix} z_1 & \cdots & z_n \end{bmatrix} \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} + b_{n1} & \cdots & a_{nn} + b_{nn} \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} =$$

$$\begin{bmatrix} z_1(a_{11} + b_{11}) + \cdots + z_n(a_{n1} + b_{n1}) & \cdots & z_1(a_{1n} + b_{1n}) + \cdots + z_n(a_{nn} + b_{nn}) \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} =$$

$$(z_1^2(a_{11} + b_{11}) + \cdots + z_1 z_n(a_{n1} + b_{n1})) + \cdots + (z_1 z_n(a_{1n} + b_{1n}) + \cdots + z_n^2(a_{nn} + b_{nn})) =$$

$$\sum_{i=1}^{n} \sum_{j=1}^{n} z_i z_j (a_{ij} + b_{ij}) =$$

$$\sum_{i=1}^{n} \sum_{j=1}^{n} z_i z_j a_{ij} + \sum_{i=1}^{n} \sum_{j=1}^{n} z_i z_j b_{ij} > 0$$

which is true because of (2) and (3). $\qquad \square$

## 5. Max Margin Hyperplanes [Exercise 7.2, page 357]

*The decision boundary $\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b = 0$ whose parameters $\boldsymbol{w}$ and $b$ are the optimal solution of the following max margin optimization problem does not depend on the value of $\gamma$.*

$$\underset{\boldsymbol{w},b}{\operatorname{argmax}} \left( \min_n \left( \frac{t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b)}{\|\boldsymbol{w}\|} \right) \right)$$
subject to
$$t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b) \geqslant \gamma$$
$$\gamma > 0.$$

---

*This box can be skipped if you prefer a shorter answer*

We assume that the training data is separable by a linear model $y(\boldsymbol{x}_n) = \boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b$ in the feature space such that $y(\boldsymbol{x}_n) \geqslant 0$ if the label for $\boldsymbol{x}_n$ is $t_n = +1$, and $y(\boldsymbol{x}_n) < 0$ if the label is $t_n = -1$. This implies that $y(\boldsymbol{x}_n) = 0$ is the decision boundary.

The *margin* is defined to be the distance of the closest point $\boldsymbol{x}_n$ to the decision surface $\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b = 0$. It can be proved that the distance of any point $\boldsymbol{x}_n$ from the surface is $\frac{|y(\boldsymbol{x}_n)|}{\|\boldsymbol{w}\|} = \frac{|\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b|}{\|\boldsymbol{w}\|}$. Now since we assume that our data is linearly separable, there exist some $\boldsymbol{w}$ and $b$ such that the model $y(\boldsymbol{x}_n) = \boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b$ can correctly classify all data points. Therefore, $t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b) > 0$ for all $\boldsymbol{x}_n$. Thus, the distance can be written as $\frac{t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b)}{\|\boldsymbol{w}\|}$ (note that $|t_n| = 1$).

To find the margin, we need to find $\boldsymbol{x}_n$ that minimizes the distance. Therefore, the margin is $\min_n \left( \frac{t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b)}{\|\boldsymbol{w}\|} \right)$ subject to $t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b) > 0$.

---

In the max margin classifier, we search for $\boldsymbol{w}$ and $b$ that maximize the margin (to minimize the generalization error), that is, $\underset{\boldsymbol{w},b}{\operatorname{argmax}} \left( \min_n \left( \frac{t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b)}{\|\boldsymbol{w}\|} \right) \right)$ subject to $t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b) > 0$.

No matter what is the optimal value for $\boldsymbol{w}$ and $b$, the scaling $\boldsymbol{w} \to \kappa\boldsymbol{w}$ and $b \to \kappa b$ does not change the classifier because it doesn't change the distances from the separating hyperplane, $\frac{t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b)}{\|\boldsymbol{w}\|}$ (because it multiplies the numerator and the denominator by the same factor $\kappa$), so it doesn't change the distance of the closest point to the hyperplane, $\min_n \left( \frac{t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b)}{\|\boldsymbol{w}\|} \right)$, either.

Therefore, we can use this scaling to set $t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b) = \gamma$ for the closest point to the hyperplane, where $\gamma$ is an arbitrary positive constant (because there exists a $\kappa$ for any given $\gamma$). That is, $\min_n \left( t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n) + b) \right) = \gamma$. This implies that (a) all points satisfy the constraint $t_n(\boldsymbol{w}^\mathsf{T}\phi(\boldsymbol{x}_n)+b) \geqslant \gamma$, and (b) to maximize the margin we need to $\underset{\boldsymbol{w},b}{\operatorname{argmax}}\frac{\gamma}{\|\boldsymbol{w}\|}$ (since $\|\boldsymbol{w}\|$ is independent of $n$) which is equivalent to $\underset{\boldsymbol{w},b}{\operatorname{argmax}}\frac{1}{\|\boldsymbol{w}\|}$ (since $\gamma$ is independent of $\boldsymbol{w}$ and $b$), or $\underset{\boldsymbol{w},b}{\operatorname{argmin}} \|\boldsymbol{w}\|$. $\qquad\square$

## 6. SVM Regression [Exercise 7.7, page 357]

SVM for regression, primal problem:

minimize $\quad J(\vec{w}, b) = \frac{1}{2}\|\vec{w}\|^2 + C\sum_{n=1}^{N}(\xi_n + \hat{\xi}_n)$

subject to $\quad t_n \leq \vec{w}^T\vec{\varphi}(\vec{x}_n) + b + \varepsilon + \xi_n \quad\Leftrightarrow\quad t_n - \vec{w}^T\vec{\varphi}(\vec{x}_n) - b - \varepsilon - \xi_n \leq 0$

$\qquad\qquad t_n \geq \vec{w}^T\vec{\varphi}(\vec{x}_n) + b - \varepsilon - \hat{\xi}_n \quad\Leftrightarrow\quad -t_n + \vec{w}^T\vec{\varphi}(\vec{x}_n) + b - \varepsilon - \hat{\xi}_n \leq 0$

$\vec{w}, b \qquad\qquad \xi_n, \hat{\xi}_n \geq 0 \qquad\qquad \forall n \in \{1, \cdots, N\} \Leftrightarrow -\xi_n, -\hat{\xi}_n \leq 0$

$$L_P(\vec{w}, \vec{\alpha}, \hat{\vec{\alpha}}, \vec{\mu}, \hat{\vec{\mu}}) = \frac{1}{2}\|\vec{w}\|^2 + C\sum_{n=1}^{N}(\xi_n + \hat{\xi}_n) + \sum_{n=1}^{N}\alpha_n(t_n - \vec{w}^T\vec{\varphi}(\vec{x}_n) - b - \varepsilon - \xi_n)$$

$$+ \sum_{n=1}^{N}\hat{\alpha}_n(-t_n + \vec{w}^T\vec{\varphi}(\vec{x}_n) + b - \varepsilon - \hat{\xi}_n) + \sum_{n=1}^{N}\mu_n(-\xi_n) + \sum_{n=1}^{N}\hat{\mu}_n(-\hat{\xi}_n)$$

derivative of $L_P$ with respect to primal variables

$\alpha_n, \hat{\alpha}_n, \mu_n, \hat{\mu}_n \geq 0 \qquad$ (Lagrange multipliers)

$\dfrac{\partial L_P}{\partial \vec{w}} = 0 \Rightarrow \vec{w} - \sum_{n=1}^{N}\alpha_n\vec{\varphi}(\vec{x}_n) + \sum_{n=1}^{N}\hat{\alpha}_n\vec{\varphi}(\vec{x}_n) = 0 \qquad$ (since $\frac{1}{2}\|\vec{w}\|^2 = \frac{1}{2}\vec{w}^T\vec{w}$,

$\qquad\qquad \Rightarrow \vec{w} = \sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n)\vec{\varphi}(\vec{x}_n) \quad$ (I) $\qquad \frac{d}{dx}(\vec{x}^T\vec{x}) = 2\vec{x}, \quad \frac{d}{dx}\vec{x}^T\vec{a} = \vec{a}$)

$\dfrac{\partial L_P}{\partial b} = 0 \Rightarrow -\sum_{n=1}^{N}(-\alpha_n + \hat{\alpha}_n) = 0 \Rightarrow \sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n) = 0 \quad$ (II)

$\dfrac{\partial L_P}{\partial \xi_n} = 0 \Rightarrow CN - \alpha_n N - \mu_n N = 0 \Rightarrow \alpha_n + \mu_n = C \quad$ (III)

$\dfrac{\partial L_P}{\partial \hat{\xi}_n} = 0 \Rightarrow CN - \hat{\alpha}_n N - \hat{\mu}_n N = 0 \Rightarrow \hat{\alpha}_n + \hat{\mu}_n = C \quad$ (IV) substitute in $L_P$ to get $L_D$

$$L_D(\vec{\alpha}, \hat{\vec{\alpha}}, \vec{\mu}, \hat{\vec{\mu}}) = \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}(\alpha_n - \hat{\alpha}_n)(\alpha_m - \hat{\alpha}_m)\vec{\varphi}(\vec{x}_n)^T\vec{\varphi}(\vec{x}_m) + \cancel{C\sum_{n=1}^{N}(\xi_n + \hat{\xi}_n)}^{\;0} \quad \text{(II)}$$

(still has $\vec{w}, b,$ $\vec{\xi},$ and $\hat{\vec{\xi}}$)

$$+ \sum_{n=1}^{N}t_n(\alpha_n - \hat{\alpha}_n) - \sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n)\vec{w}^T\vec{\varphi}(\vec{x}_n) - \underbrace{b\sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n)}_{\text{(III)}}$$

$$- \varepsilon\sum_{n=1}^{N}(\alpha_n + \hat{\alpha}_n) - \underbrace{\sum_{n=1}^{N}(\alpha_n + \mu_n)}_{C \;\text{(III)}}\xi_n - \underbrace{\sum_{n=1}^{N}(\hat{\alpha}_n + \hat{\mu}_n)}_{C \;\text{(IV)}}\hat{\xi}_n$$

$$= \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}(\alpha_n - \hat{\alpha}_n)(\alpha_m - \hat{\alpha}_m)\vec{\varphi}(\vec{x}_n)^T\vec{\varphi}(\vec{x}_m) + \sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n)t_n$$

$$- \sum_{n=1}^{N}\sum_{m=1}^{N}(\alpha_n - \hat{\alpha}_n)(\alpha_m - \hat{\alpha}_m)\vec{\varphi}(\vec{x}_n)^T\vec{\varphi}(\vec{x}_m) - \varepsilon\sum_{n=1}^{N}(\alpha_n + \hat{\alpha}_n)$$

Dual Problem

maximize $\quad L_D(\vec{\alpha}, \hat{\vec{\alpha}}) = -\frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N}(\alpha_n - \hat{\alpha}_n)(\alpha_m - \hat{\alpha}_m)k(\vec{x}_n, \vec{x}_m) \quad$ (**)

(**) $k(\vec{x}_n, \vec{x}_m) = \vec{\varphi}(\vec{x}_n)^T\vec{\varphi}(\vec{x}_m)$

$\qquad\qquad\qquad\qquad - \varepsilon\sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n) + \sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n)t_n$

— (I)

$y(\vec{x}) = \sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n)k(\vec{x}, \vec{x}_n) + b \quad$ (III)(IV)

subject to

$0 \leq \alpha_n \leq C$

$0 \leq \hat{\alpha}_n \leq C$

$\sum_{n=1}^{N}(\alpha_n - \hat{\alpha}_n) = 0$

$y(\vec{x}) = \vec{w}\vec{\varphi}(\vec{x}) + b \quad$ (II)

Complementary slackness

$\alpha_n(t_n - y_n - \varepsilon - \xi_n) = 0$

$\hat{\alpha}_n(-t_n + y_n - \varepsilon - \hat{\xi}_n) = 0$

$(C - \alpha_n)\xi_n = 0 \quad$ — (III)

$(C - \hat{\alpha}_n)\hat{\xi}_n = 0 \quad$ — (IV)

(*)

## 7. Large Margin Perceptron

*Let $\mathbf{u}$ be a current vector of parameters and $\mathbf{x}$ and $\mathbf{y}$ two training examples such that $\mathbf{u}^T(\mathbf{x} - \mathbf{y}) < 1$ . Use the technique of Lagrange multipliers to find a new vector of parameters $\mathbf{w}$ as the solution to the convex optimization problem below: minimize: $J(\mathbf{w}) = \frac{1}{2}\|\mathbf{w} - \mathbf{u}\|^2$ subject to: $\mathbf{w}^T(\mathbf{x} - \mathbf{y}) \geqslant 1$*

(⊛)

$\vec{u}$ : current vector of parameters

$\vec{x}, \vec{y}$ : two training examples $\qquad \vec{u}^T(\vec{x} - \vec{y}) < 1 \qquad$ (⊛⊛)

$\vec{w}$ : new vector of parameters $\qquad \longleftarrow$ primal problem

minimize $\quad J(\vec{w}) = \frac{1}{2}\|\vec{w} - \vec{u}\|^2$

subject to $\quad \vec{w}^T(\vec{x} - \vec{y}) \geqslant 1 \qquad \Rightarrow \quad 1 - \vec{w}^T(\vec{x} - \vec{y}) \leqslant 0$

$L_P(\vec{w}, \alpha) = \frac{1}{2}\|\vec{w} - \vec{u}\|^2 + \alpha\left(1 - \vec{w}^T(\vec{x} - \vec{y})\right) \qquad \longleftarrow$ primal Lagrangian

$\qquad = \frac{1}{2}(\vec{w} - \vec{u})^T(\vec{w} - \vec{u}) + \alpha - \alpha \vec{w}^T(\vec{x} - \vec{y}) \qquad \vec{u}, \vec{x}, \vec{y}$ : constant vectors

$\underbrace{\qquad}_{\text{chain rule}}$

$\frac{\partial L_P}{\partial \vec{w}} = (\frac{1}{2})(2(\vec{w} - \vec{u})\mathbb{1}) - \alpha(\vec{x} - \vec{y}) = 0 \qquad \frac{d}{\vec{x}}\vec{x}^T\vec{x} = 2\vec{x}$ | set the derivative of Lagrangian

$\Rightarrow \quad \vec{w} - \vec{u} = \alpha(\vec{x} - \vec{y}) \qquad\qquad \frac{d}{d\vec{x}}\vec{x}^T\vec{a} = \vec{a}$ | with respect to the primal variable

$\Rightarrow \quad \vec{w} = \vec{u} + \alpha(\vec{x} - \vec{y}) \quad$ (⊛⊛⊛) | to zero

$L_D(\alpha) = \frac{\alpha^2}{2}\|\vec{x} - \vec{y}\|^2 + \alpha - \alpha\left(\vec{u} + \alpha(\vec{x} - \vec{y})\right)^T(\vec{x} - \vec{y})$ | substitute primal variable back in $L_P$ to get $L_D$

$\qquad = \frac{\alpha^2}{2}\|\vec{x} - \vec{y}\|^2 + \alpha - \alpha\vec{u}^T(\vec{x} - \vec{y}) - \alpha^2\|\vec{x} - \vec{y}\|^2$

$\qquad = -\frac{\|\vec{x} - \vec{y}\|^2}{2}\alpha^2 + \left(1 - \vec{u}^T(\vec{x} - \vec{y})\right)\alpha$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \longleftarrow$ dual problem

maximize $\quad L_D(\alpha) = -\frac{\|\vec{x} - \vec{y}\|^2}{2}\alpha^2 + \left(1 - \vec{u}^T(\vec{x} - \vec{y})\right)\alpha$

subject to $\qquad \alpha \geqslant 0$ | a maximization problem in one variable.

$\frac{dL_D}{d\alpha} = 0 \quad \Rightarrow \quad -\|\vec{x} - \vec{y}\|^2\alpha + 1 - \vec{u}^T(\vec{x} - \vec{y}) = 0$ | set the derivative to zero to compute optimal $\alpha$ | we can check that objective function is concave

$\qquad\qquad \Rightarrow \quad \alpha = \frac{1 - \vec{u}^T(\vec{x} - \vec{y})}{\|\vec{x} - \vec{y}\|^2}$ | $\frac{d^2 L_D}{d\alpha^2} = -\|\vec{x} - \vec{y}\|^2 < 0$

$\alpha \geqslant 0$ because $\|\vec{x} - \vec{y}\|^2 > 0 \quad$ (*)

$\qquad$ and $1 - \vec{u}^T(\vec{x} - \vec{y}) > 0 \quad$ (⊛⊛)

$\boxed{\vec{w} = \vec{u} + \frac{1 - \vec{u}^T(\vec{x} - \vec{y})}{\|\vec{x} - \vec{y}\|^2}(\vec{x} - \vec{y})}$ | substitute $\alpha$ back into (⊛⊛⊛)

## 8. Feature Selection (Document Classifier)

(a) The data is obtained from UC Irvine machine learning, Dexter dataset.

The data is derived from a text corpus, formatted in the bag-of-word representation. There are a total of $20,000$ features, where each feature represents the frequency of a word stem used in text.

Only the training and validation data are useful to us because the labels are provided for them. Each of these two sets includes 300 examples where half of them are labeled 1 (meaning the text is related to "corporate acquisition") and the other half -1.

There is a considerable number of features that are zero for all examples. Around half of the features are real text attributes, and the other half are just dummy features that were added to the data in a randomized manner to make for a feature selection exercise.

(b) According to the dataset description, the features represent "frequencies of occurrences of word stems in text". So the feature values are relevant as they indicate some sort of measurement, and thus they are not nominal features.

(Nominal numbers uniquely identify an object or entity. For instance, zip codes, or social security numbers identify geographic areas or individuals, respectively. Nominal numbers can be compared only for equality–to see if two geographic areas or two people are actually the same–but the value of the numbers do not mean anything. So comparison of their relative value, their difference, their sum, etc. is nonsensical.)

Therefore, we cannot use Mutual Information and Chi-square Statistic which can handle only nominal features, but we can use Pearson Correlation Coefficient, Signal-to-Noise Ratio, and T-test.

We use only the training data for ranking features. So we rank features using each of the three measures of correlation with the label. Then we select the top N features to be used for training linear SVM and the nearest neighbor algorithm.

The training data was formatted into the file data/svm-data.train using the Python script code/preprocess-data-libsvm.py to fit the format required by LIBSVM, then it was read into a sparse matrix using the Python script code/rank-features.py.

The feature matrix for each data set (that contains 300 examples of each $20,000$ features) has a total of $6,000,000$ elements of which only a small minority of $28,218$ elements are not zero. Therefore, storing the feature matrix as a sparse matrix seems a natural choice. This was done using

the `scipy.sparse` library.

**Pearson Correlation Coefficient**

The mean value for labels, $\overline{y}$, of Dexter dataset is 0 as half of the examples are labeled 1 and the other half are labeled -1.

The absolute value of the correlation coefficient matters to us because we only care about the magnitude of correlation between a feature and the label. For the purpose of feature selection, we do not care if the correlation is direct or inverse.

If the correlation coefficient is undefined (if the denominator is 0), we replace the associated feature with 0 which implies the label is independent of that feature.

The division by zero in the correlation coefficient

$$\rho(X, Y) = \frac{\sum_{i=1}^{n}(X_i - \overline{X})(Y_i - \overline{Y})}{\sqrt{\sum_{i=1}^{n}(X_i - \overline{X})^2}\sqrt{\sum_{i=1}^{n}(Y_i - \overline{Y})^2}}$$

happens if and only if $\sum_{i=1}^{n}(X_i - \overline{X})^2 = 0$ because $\sum_{i=1}^{n}(Y_i - \overline{Y})^2$ cannot be zero for all examples of a classification problem with at least two classes. Now, $\sum_{i=1}^{n}(X_i - \overline{X})^2 = 0$ happens if and only if all nonnegative $(X_i - \overline{X})^2$ terms are zero (for all $i = 1, \ldots, n$). In which case, the numerator has to be zero, too. So all division by zero cases are actually a $\frac{0}{0}$ case for the correlation coefficient. And this happens for a feature that has the same value across all examples. These features are useless for training purposes. So by manually setting their correlation coefficient to zero, they will be ranked last for feature selection.

The absolute value of the correlation coefficient is between 0 and 1.

**Signal-to-Noise Ratio**

In the `rank-features.py` script, we use the following formula to compute the standard deviation of a feature X.

$$\sigma = \sqrt{\overline{X^2} - \overline{X}^2}$$

Signal-to-Noise (S2N) ratio of feature X in relation to label Y is then computed as

$$\mu(X, Y) = \frac{|\mu_+ - \mu_-|}{\sigma_+ + \sigma_-}$$

where $\mu_+$ and $\sigma_+$ are the sample mean and standard deviation of X for which $Y = 1$.

All S2N ratio values for our data is between 0 and 1.

### T-test

$$T(X, Y) = \frac{|\mu_+ - \mu_-|}{\sqrt{\sigma_+^2/m_+ + \sigma_-^2/m_-}}$$

where $m_+$ is the number of examples with label 1.

The maximum value of t-test can be larger than 1.

### Sort columns of feature matrix

We sort the columns of the feature matrix to have the most relevant feature on the first column of the matrix to the least relevant on the last column.

This sorting is done based on each of the three filters. The sort will be based on the value of the filter in descending order. That is, the feature with the highest filter value will be moved to the left-most column in the matrix.

### Normalization

After feature selection, we divide each feature vector $x_i$ by its Euclidean norm, so all of the resulting feature vectors $\frac{x_i}{|x_i|}$ will have a magnitude of 1.

### Sort features first, then normalize.

We first sort the column indices of the feature matrix using a filter as sorting keys (a filter value will be computed for each column). We then normalize features vectors (matrix rows) in the feature matrix, and finally sort the columns of the normalized feature matrix by applying the sorted column indices.

This sequence of operations allows us to carry out the normalization only once regardless of the number of filters that we wish to use.