Dhruv Chakraborty
UID: 204-962-098
Prof. Sankararaman

CM146: PSET 3

1) <u>Kernels</u>

(a) $k(n,z)$ is a kernel. We can see that it is symmetric since $k(n,z) = k(z,n)$ as the intersection of the words remains the same. Now we show positive semi-definiteness.

set $K = \begin{bmatrix} k(n,n) & k(n,z) \\ k(z,n) & k(z,z) \end{bmatrix}$  ∵ $k(n,z) = k(z,n)$

$$\det(K - \lambda I) = (k(n,n) - \lambda)(k(z,z) - \lambda) - k(n,z)^2 = 0$$
$$\Rightarrow \lambda^2 - \lambda(k(n,n) + k(z,z)) + k(n,n)k(z,z) - k(n,z)^2 = 0$$

Using quadratic formula to solve for $\lambda$,

$$\lambda = \frac{k(n,n) + k(z,z) \pm \sqrt{(k(n,n) + k(z,z))^2 - 4(k(n,n)k(z,z) - k(n,z)^2)}}{2}$$

$$= \frac{k(n,n) + k(z,z) \pm \sqrt{}}{}$$

Clearly, $k(n,n) + k(z,z)$ or $\sqrt{(k(n,n) + k(z,z))^2 - 4(k(n,n)k(z,z) - k(n,z)^2)}$

is positive so $\lambda > 0$ (as $k$ is +ve and $\sqrt{n} > 0 \, \forall n \in \mathbb{R}$)

Now considering $\lambda_2 = \frac{k(n,n) + k(z,z) - \sqrt{(k(n,n) + k(z,z))^2 - 4(k(n,n)k(z,z) - k(n,z)^2)}}{2}$

We want $k(n,n) + k(z,z) \geq \sqrt{(k(n,n) + k(z,z))^2 - 4(k(n,n)k(z,z) - k(n,z)^2)}$

⇕

$(k(n,n) + k(z,z))^2 \geq (k(n,n) + k(z,z))^2 - 4(k(n,n)k(z,z) - k(n,z)^2)$

But ⇕

$0 \geq -4(k(n,n)k(z,z) - k(n,z)^2)$

⇕

$k(n,n)k(z,z) \geq k(n,z)^2$

But we know this since the intersection of a document with itself is necessarily greater than/equal to a different document i.e. $k(n,n) \geq k(n,z)$ and $k(z,z) \geq k(n,z)$.

Since $K$ is symmetric and has only ~~+ve~~ non-ve eigenvalues it is positive semi-definite.

∴ $k(n,z)$ is a kernel.

(b) Since $k(n,z) = n \cdot z$ is a kernel, we use scaling to with $f(n) = \dfrac{1}{\|n\|}$ to get that $\dfrac{1}{\|n\|} \, n \cdot z \, \dfrac{1}{\|z\|}$ is also a kernel.

Also since $k^*(n,z) = 1$ is clearly a kernel, $1 + \dfrac{n \cdot z}{\|n\| \, \|z\|}$ is a kernel through the sum rule.

Multiplying this kernel by itself thrice and using the product rule, $\left(1 + \left(\dfrac{n}{\|n\|}\right) \cdot \left(\dfrac{z}{\|z\|}\right)\right)^3$ is a kernel.

---

(c) $k_\beta(n,z) = (1 + \beta \vec{n} \cdot \vec{z})^3$

$= (1 + \beta(n_1 z_1 + n_2 z_2))^3 \quad \because n, z \in \mathbb{R}^2$

$= (1 + \beta n_1 z_1 + \beta n_2 z_2)^2 (1 + \beta n_1 z_1 + \beta n_2 z_2)$

$= (1 + \beta^2 n_1^2 z_1^2 + \beta^2 n_2^2 z_2^2 + 2\beta n_1 z_1 + 2\beta n_2 z_2 + 2\beta^2 n_1 z_1 n_2 z_2)$
$\quad (1 + \beta n_1 z_1 + \beta n_2 z_2)$

$= 1 + \beta^2 n_1^2 z_1^2 + \beta^2 n_2^2 z_2^2 + 2\beta n_1 z_1 + 2\beta n_2 z_2 + 2\beta^2 n_1 z_1 n_2 z_2 +$
$\beta n_1 z_1 + \beta^3 n_1^3 z_1^3 + \beta^3 n_1 n_2^2 z_1 z_2^2 + 2\beta^2 n_1^2 z_1^2 + 2\beta^2 n_1 n_2 z_1 z_2 +$
$2\beta^3 n_1^2 z_1^2 n_2 z_2 + \beta n_2 z_2 + \beta^3 n_1^2 z_1^2 n_2 z_2 + \beta^3 n_2^3 z_2^3 +$
$2\beta^3 n_1 z_1 n_2 z_2 + 2\beta^3 n_2^2 z_2^2 + 2\beta^3 n_1 z_1 n_2^2 z_2^2$

$= 1 + 3\beta n_1 z_1 + 3\beta^2 n_1^2 z_1^2 + \beta^3 n_1^3 z_1^3 + 3\beta n_2 z_2 + 6\beta^2 n_1 z_1 n_2 z_2$
$+ 3\beta^3 n_1^2 z_1^2 n_2 z_2 + 3\beta^2 n_2^2 z_2^2 + 2\beta^3 n_1 z_1 n_2^2 z_2^2 + \beta^3 n_2^3 z_2^3$

So $\varphi_\beta(n) =$

$$
\begin{bmatrix}
1 \\
\sqrt{3\beta}\, n_1 \\
\sqrt{3}\beta\, n_1^2 \\
\beta^{3/2}\, n_1^3 \\
\sqrt{3\beta}\, n_2 \\
\sqrt{6}\,\beta\, n_1 n_2 \\
\sqrt{3}\,\beta^{3/2}\, n_1^2 n_2 \\
\sqrt{3}\,\beta\, n_2^2 \\
\sqrt{3}\,\beta^{3/2}\, n_1 n_2^2 \\
\beta^{3/2}\, n_2^3
\end{bmatrix}
$$

→ The difference between $k_\beta(n,z)$ and $k(n,z)$ is that $k_\beta$ has the scaling factor $\beta$ that scales the vector. It could potentially be used to regularize the ~~tool~~

$\Rightarrow$ $\beta$ scales the vector.

2) <u>SVM</u>

(a) Want to minimize $\frac{1}{2}\|\theta\|^2$ with $y_n\theta^T x_n \ge 1$ i.e. $-\theta^T(a,e)^T \ge 1$

$\Rightarrow 1 + a\theta_1 + e\theta_2 \le 0$.

$L(\theta, \alpha) = \frac{1}{2}\|\theta\|^2 + \alpha(\theta (a,e)^T + 1)$

So $\frac{\partial L}{\partial \theta} = \theta + \alpha(a,e)^T = 0 \Rightarrow \theta = -\alpha\begin{pmatrix} a \\ e \end{pmatrix}$

Now maximizing over $\alpha$:

$\max_\alpha \frac{1}{2}\left\| -\alpha\begin{pmatrix} a \\ e \end{pmatrix}\right\|^2 + \alpha\left(-\alpha\begin{pmatrix} a \\ e \end{pmatrix}(a\ e) + 1\right)$

$= \max_\alpha \frac{1}{2}(\alpha^2 a^2 + \alpha^2 e^2) - (\alpha^2 a^2 + \alpha^2 e^2) + \alpha$

$= \max_\alpha -\frac{\alpha^2(a^2 + e^2)}{2} + \alpha$

$\frac{\partial}{\partial \alpha} = -\alpha(a^2 + e^2) + 1 = 0 \Rightarrow \alpha^* = \frac{1}{a^2 + e^2}$

So $\theta^* = \frac{-1}{a^2 + e^2}\begin{pmatrix} a \\ e \end{pmatrix}$

(b) Given $x_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, $x_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $y_1 = 1$ and $y_2 = -1$ we must

satisfy $\theta_1 + \theta_2 \ge 0$, $\theta_1 \le 0$ and

$\theta_1 + \theta_2 \ge 1$, $\theta_1 \le 1$ for $y_n \theta^T x_n \ge 1$

$\theta_1 = -1$ and $\theta_2 = 2$ satisfy these constraints

$\therefore \theta^* = \begin{pmatrix} -1 \\ 2 \end{pmatrix}$ and $\gamma = \frac{1}{\sqrt{1^2 + 2^2}} = \frac{1}{\sqrt{5}}$

(c) The constraints with the offset term are $\theta_1 + \theta_2 \ge 1 - b$ & $\theta_1 \le 1 - b$

These are satisfied by $\theta = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$ and $b = -1$.

The margin $\gamma = \frac{1}{\sqrt{2^2}} = \frac{1}{2}$, which is larger than the margin

without the offset $\left(\frac{1}{\sqrt{5}}\right)$

2) Twitter Analysis using SVMs

3.2) Hyperparameter Selection for a Linear-Kernel SVM

(b) Since we wish to find hyperparameters that generalise best to unseen data, specifically measured using training/test data, we want the dataset we are testing on to have a similar distribution. A fold that doesn't have the same proportions could lead to overfitting and is an would be an outlier. Moreover, a fold with mostly +ve/-ve samples teaches the model nothing.

(d)

| C | accuracy | F1-score | AUROC | precision | sensitivity | specificity |
|---|---|---|---|---|---|---|
| $10^{-3}$ | 70.89% | 82.97% | 50% | 70.89% | 100% | 0% |
| $10^{-2}$ | 71.07% | 83.06% | 50.31% | 71.02% | 100% | 0.63% |
| $10^{-1}$ | 80.60% | 87.55% | 71.88% | 83.57% | 92.94% | 50.81% |
| $10^{0}$ | 81.46% | 87.49% | 75.31% | 85.62% | 90.17% | 60.45% |
| $10^{1}$ | 81.82% | 87.66% | 75.92% | 85.95% | 90.17% | 61.67% |
| $10^{2}$ | 81.82% | 87.66% | 75.92% | 85.95% | 90.17% | 61.67% |
| best C | 10 | 100 | 10 | 10 | 0.001 | 10 |

On all performance metrics $C=10$, ~~100~~ and $C=100$ give the same value, ~~best~~ which is actually the best for all metrics outside sensitivity. For sensitivity, smaller values like $C=10^{-3}$ and $c=10^{-2}$ perform better, but this is an ~~ett~~ ~~major~~ outlier and susceptible to overfitting. We use $C=10$ as our ~~best~~ best C.

3.3) Hyperparameter Selection for an RBF-Kernel SVM

(a) The $\gamma$ parameter "defines how far the influence of a single training example reaches, lower values meaning far and higher values meaning close." (from sklearn docs). Basically as $\gamma$ increases the each support vectors use radius of influence decreases. We can tune it to avoid overfitting and increase generalization.

(b) In my grid, both C and γ range from $10^{-3}$ to $10^2$ in powers of 10. We used these values for C for the linear kernel and they performed quite well. Moreover, this grid allows for a very large range of values for both parameters.

(c)

| metric | score | C | γ |
|---|---|---|---|
| accuracy | 79.34% | 100 | 0.01 |
| F1-score | 87.63% | 100 | 0.01 |
| AUROC | 75.45% | 100 | 0.01 |
| precision | 85.83% | 100 | 0.01 |
| sensitivity | 100% | 0.001 | 0.001 |
| specificity | 60.47% | 100 | 0.01 |

Just as in the linear kernel, sensitivity is an outlier with C=0.001 and γ=0.001 performing best for it. However, for every other performance metric, C=100 and γ=0.01 are the best hyperparameters in our grid.

3.4) Test Set Performance

(a) We use C=10 for the linear kernel SVM as it performed the best on all metrics except sensitivity (as did C=100). Similarly, as C=100 & γ=0.01 did the best for almost all metrics for the RBF kernel SVM, we use these hyperparameters to train it.

(c)

| metric | linear | RBF | |
|---|---|---|---|
| accuracy | 74.29% | 75.71% | As we can see, the RBF |
| F1-score | 43.75% | 45.16% | kernel classifier performs |
| AUROC | 62.59% | 63.61% | slightly better than the |
| precision | 63.64% | 70.00% | linear kernel SVM on |
| sensitivity | 33.33% | 33.33% | all metrics (except sensitivity). |
| specificity | 91.84% | 93.88% | We would ideally deploy |
| | | | this model over the linear one |