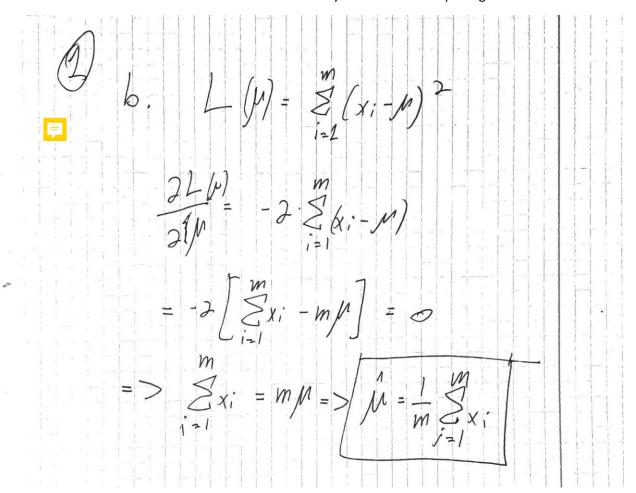
- 1.
- a. K-medoid is more robust to noise and outliers because it uses an actual data point as the cluster center rather than the cluster mean which may be very far away from the actual cluster center when using data with large outliers, furthermore K-medoid allows for using other evaluation metrices rather than the Euclidean distance, which maybe better for clustering (for example using some feature correlation between 2 data points).
- b. We can find the minimum of the term by derivation and equating to zero:



Bouns: *** attached at the end of document ***

- A-1
- D-2
- C 3
- E-4
- B-5
- F 6

<u>Explanation</u>: A & D are the linear kernels therefore correspond to labels 1 & 2, we can see that in A the decision boundary passes through 2 of the data points and in D it doesn't touch any data point, therefore we understand that A's SVM had more "slack" in its training meaning lower higher impact of ξ therefore lower C.

B & F are the RBF kernels because of the decision boundary shape which can only be achieved using radial functions, F is overfitted thus has the higher Gamma.

C & E are the polynomial kernels, C is the 2nd order because its decision boundary is a parabola, therefore E is the 10th must be the 10th order.

3.

- a. The scientific term referred to is bias-variance tradeoff.
- =
- b. The ln(L) term correlates to the bias as it states how well does our data fit with our model given a set of chosen parameters, the P term refers to the variance because it counts the number of parameters in our model.
- c. Violation of the bias-variance tradeoff can result in either overfitting or underfitting.
- d. We want to minimize the AIC because we want best fit for our model (meaning high ln(L)) while having the simplest model – which means fewer parameters as possible – meaning low P.

$$\frac{1}{1} \left(\frac{1}{1} x - \frac{1}{1} \right)$$

If we take the sirst and last terms out of the sum we get.

$$L(n) = \sum_{i=2}^{m-1} (|x_i, n|) + |x_i - n| + |x_i - n|$$

$$= \sum_{i=2}^{m-1} \{ [x, -n] \} + (x_n - n) - (x, -n)$$

$$= \underbrace{\begin{array}{c} m^{-1} \\ X_{1} - m \end{array}}_{i=2} + \underbrace{\left(X_{n} - X_{i} \right)}_{i=2}$$

we repeat this process until he have for 2 elements lest in the sum (depending

is m is odd then we get:

In this case the minimum is given by N=Xn+1 which is the group median Is mis evenue get: L(M) = |Xn - M | + | Xn+2 - M + const. In this case the minimum is received Sor the mean of the two sumples. $\hat{\Lambda} = \frac{x_{1}^{m} + x_{1}^{m+2}}{2}$ which is the median aswell Finally we get: $\mathcal{N} = \begin{cases} \times \frac{m+1}{2} \\ \times \frac{m}{2} + \times \frac{m}{2} \end{cases}$ m isodd m is even