

To find the posterior density $p(\lambda|\mathcal{X})$ we need a prior on λ . We claim that a conjugate prior for the exponential distribution is the gamma distribution

$$\text{Gamma}(\lambda|\alpha, \beta) = \frac{\Gamma(\alpha)}{\beta^\alpha} \lambda^{\alpha-1} \exp(-\lambda\beta),$$

where $\Gamma(\alpha) = \int_0^\infty \exp^{-t} t^{\alpha-1} dt$ is the gamma function.

- b) 1. What does *conjugate prior* mean? 1 pt.

it means the posterior would have the same form
as the prior

2. Show that the gamma distribution is the conjugate prior of
the exponential distribution. 2 pts.

$$p(\lambda|x) \propto p(x|\lambda) p(\lambda|\beta, \gamma) \propto \lambda^{\alpha-1} \cdot \lambda^{\beta-1} \cdot \exp^{-\lambda(\beta + x_1 + \dots + x_n)}$$

which is also a gamma function with $\alpha' = \alpha + n$
 $\beta' = \beta + x_1 + \dots + x_n$

- c) Given a Gamma prior over the rate λ (prior with parameters α and β), write the maximum a posteriori estimator $\hat{\lambda}_{\text{MAP}}(\mathcal{X})$ as an explicit function of the i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$:
(please write the direct closed-form solution.)

$$\hat{\lambda}_{\text{MAP}}(\mathcal{X}) = \arg \max p(\lambda|\mathcal{X}) = \text{2 pts.}$$

$$\arg p(\lambda|x) \propto (\alpha+n-1) \log \lambda - \lambda(\beta+x_1+\dots+x_n)$$

$$\frac{\frac{\partial \log p(\lambda|x)}{\partial \lambda}}{\frac{\partial}{\partial \lambda}} = \frac{\alpha+n-1}{\lambda} - (\beta+x_1+\dots+x_n) = 0$$

$$\Rightarrow \lambda = \frac{\alpha+n-1}{\beta+x_1+\dots+x_n}$$

- d) When is the maximum likelihood estimator (MLE) equal to the maximum a posteriori (MAP) estimator given a set of i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$?

1. If the number of observations is finite.

2 pts.

$$\frac{n}{x_1 + \dots + x_n} = \frac{\alpha n - 1}{\beta + x_1 + \dots + x_n} \quad \text{when this equation is right!}$$

2. If the number of observations is infinite ($n \rightarrow \infty$). 2 pts.

$$\lim_{n \rightarrow \infty} \frac{n}{x_1 + \dots + x_n} = \frac{\alpha n - 1}{\beta + x_1 + \dots + x_n}$$

in this case ~~$\lambda = \beta$~~

see according to the convergence of MLE,

$$\lim_{n \rightarrow \infty} \frac{n}{x_1 + \dots + x_n} = \lambda$$

in this case it will automatically

- e) Assume that you have a set of the i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$ and that you can not decide which distribution to use for data description: the Gaussian distribution or the Beta distribution.

If you use the Bayesian framework what you can look at? 2 pts.

$$P(\lambda | \mathcal{X}) = \frac{P(x|\lambda) \pi(\lambda)}{\text{prior}}$$

compute the likelihood function and the prior value for the MLE

Now consider a binary classification task from a set of the i.i.d. observations $\mathcal{X} = \{x_i, y_i\}_{i=1}^n$, with $x \in \mathbb{R}^D$. Assume that the likelihood of both classes is Gaussian (assume class prior π_i , mean μ_i , and covariance matrix Σ_i for class y_i , with $i = 1, 2$).

- f) Recall that a discriminant function for class y_i is defined as:

$$g_{y_i}(x) = p(y_i|x).$$

How can you find a decision surface in terms of likelihood, prior and evidence? 1 pt.

$$p(y_1|x) = p(y_2|x)$$

$$\underbrace{p(x|y_1) p(y_1)}_{\text{prior}} = \frac{p(x|y_1) p(y_1)}{4}$$

$$p(x|y_1) \pi_1 = p(x|y_2) \pi_2$$

$$\frac{1}{2\sigma^2} \|x - \mu_1\|^2 +$$

g) Assume that

$$\mu_1 = \mu_2 = \mu$$

$$\Sigma_1 = \frac{1}{2\lambda_1} \mathbb{I}, \quad \Sigma_2 = \frac{1}{2\lambda_2} \mathbb{I}$$

$$\lambda_1 > 0, \quad \lambda_2 > 0, \quad \lambda_1 \neq \lambda_2$$

where \mathbb{I} denotes the identity matrix. Write the equation satisfied by the separating decision surface. The equation must be an explicit function of x_1 (the single observation), of the class prior, means and covariance:

(please write the solution in the polynomial form)

1. Decision surface:

3 pts.

$$g_{y_1}(x) = \log p_{y_1}(x) + \log \pi_{y_1} = -\frac{1}{2} \lambda_1 \|x - \mu\|^2 + \log \pi_{y_1} + \frac{d}{2} \log \lambda_1$$

$$g_{y_1}(x) = g_{y_2}(x)$$

$$\Rightarrow \lambda_1 \|x\|^2 + 2u_{y_1} \cdot x - \|u_{y_1}\|^2 + \log \pi_{y_1} + \frac{d}{2} \log \lambda_1 = \lambda_2 \|x\|^2 + 2u_{y_2} \cdot x - \|u_{y_2}\|^2 + \log \pi_{y_2} + \frac{d}{2} \log \lambda_2$$

$$\Rightarrow x_1 = \pm \sqrt{\frac{\log \frac{\pi_{y_1}}{\pi_{y_2}} + \frac{d}{2} \log \frac{\lambda_2}{\lambda_1}}{\lambda_1 - \lambda_2}} \quad (x - u)^T = \frac{-\log \frac{\pi_{y_1}}{\pi_{y_2}}}{\lambda_1 - \lambda_2}$$

2. In the case described above, is the decision surface linear, parabolic, spherical, cylindrical, or something else?

linear | parabolic | spherical | cylindrical | other

2 pts.

linear

depend on the covariance matrix

and prior

Question 2: Linear Classifiers and Kernels (20 pts.)

- a) Below is a list of algorithms which given a training set output a prediction function. Cross **all** of the algorithms that necessarily output a linear (in the original space) prediction function.

- Neural network
- Perceptron with learning rate $\eta = 1$
- SVM with radial basis kernel
- K-nearest neighbor classifier
- SVM with polynomial kernel with degree 1
- Ridge regression

3 pts.

- b) Recall the SVM problem. As a constrained optimization problem a solution can be obtained through both the primal and the dual form.

1. Given a primal solution for the SVM, write down the resulting classifier. 1 pt.

$$g(x) = \sum_{i=1}^n z_i z_i x^T x + w_0$$

$$w_0 = -\frac{1}{2} \left[\min_{z_i=1} w^T y_i + \max_{z_i=-1} w^T y_i \right]$$

$$\sum z_i z_i = 0$$

$$y = \underbrace{\text{sign}(w^T x + w_0)}$$

2. Given a dual solution for the SVM, write down the resulting classifier. 1 pt.

$$g(x) = \underbrace{\sum_{i=1}^n z_i z_i x^T x + w_0}$$

$$\sum z_i z_i = 0$$

z_i^* is the optimal solution

$$w = \sum_{i=1}^n z_i^* z_i x_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n z_i z_j z_j^* y_i^T y_j$$

$$z_i > 0$$

3. In practice, often the dual form of the SVM is solved to obtain a classifier. Provide one advantage of solving the dual SVM instead of the primal. 3 pts.

Since the solution to dual SVM induces sparsity to the decision function, where only few support vector enter the sum to evaluate the function

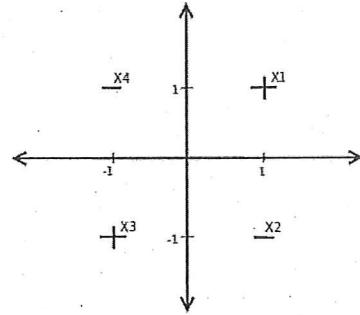
- c) Let $S = \{\mathbf{x}_i, y_i\}_{i=1}^4$ be the following training set - efficiency and apply kernel tricks interpretability

$$\mathbf{x}_1 = (1, 1) \quad y_1 = 1,$$

$$\mathbf{x}_2 = (1, -1) \quad y_2 = -1,$$

$$\mathbf{x}_3 = (-1, -1) \quad y_3 = 1,$$

$$\mathbf{x}_4 = (-1, 1) \quad y_4 = -1$$



Suppose that we trained an SVM on S and the resulting classifier $f(x)$ achieved zero training error. We ask you to provide an explicit description of $f(x)$ (a formula with numeric values).

Hint: Think of a suitable kernel function or alternatively a feature map. 5 pts.

$$(\mathbf{x}_1, \mathbf{x}_2) \rightarrow \mathbf{x}_1 \cdot \mathbf{x}_2$$

$$\mathbf{x}_1 \mathbf{x}_1' \mathbf{x}_2 \mathbf{x}_2'$$

$$(z_1 z_2 \phi(\mathbf{x})^T \phi(\mathbf{x}) + w_0)$$

$$k = \mathbf{x}_1' \mathbf{x}_1 + 9 \mathbf{x}_1' \mathbf{x}_2 \mathbf{x}_2' + \mathbf{x}_2' \mathbf{x}_2 = \min_{z_1, z_2} (w^T \phi(\mathbf{x}_1) + w^T \phi(\mathbf{x}_2))$$

$$\phi(\mathbf{x}) = [\mathbf{x}_1, \mathbf{x}_1 \mathbf{x}_2, \mathbf{x}_2]^T$$

$$\boxed{\mathbf{x}_1 + \mathbf{x}_2 + 3\mathbf{x}_1 \mathbf{x}_2}$$

$$\phi(\mathbf{x}) = (1, 1, 1)^T$$

$$z_1 z_2 = 0$$

$$\begin{aligned} w^T \phi(\mathbf{x}_1) &= -\frac{1}{2} (1, 1, 1)^T (1, 1, 1) + 1(1, 1, 1)^T (1, 1, 1) \\ &= -\frac{1}{2} (2+2+1) + 1 = 0 \end{aligned}$$

$$z_1 + z_2 = z_2 + z_4$$

$$\text{Let } z_1 = x_2 = 1$$

$$\begin{aligned} w^* &= \sum_{i=1}^4 z_i z_i \phi(\mathbf{x}_i)^T \\ &= (1, 1, 1)^T + (-1, 1, -1)^T + (-1, -1, 1)^T + (-1, -1, -1)^T \\ &= (1, 0, 2)^T - (1, -1, 1)^T \\ &= (1, 1, 1)^T - (1, -1, 1)^T \\ &= (2, 2, 0)^T \end{aligned}$$

$$\text{or } (\mathbf{x}_1, 3\mathbf{x}_1 \mathbf{x}_2)$$

$$k = \mathbf{x}_1' \mathbf{x}_1 + 9 \mathbf{x}_1' \mathbf{x}_2 \mathbf{x}_2'$$

$$\begin{aligned} \phi(\mathbf{x}_1) &= (1, 1, 1)^T \\ \phi(\mathbf{x}_2) &= (1, -1, -1)^T \\ \phi(\mathbf{x}_3) &= (-1, 1, -1)^T \\ \phi(\mathbf{x}_4) &= (-1, -1, 1)^T \end{aligned}$$

$$\begin{aligned} k &= (1, 1, 1)^T + 9(1, 1, 1)^T (1, -1, -1)^T + 2(-1, 1, -1)^T (-1, -1, 1)^T \\ &= 10(-1, 3, 3) + (-1, 3, -3) \end{aligned}$$

d) Consider a training set $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \{-1, 1\}$

1. Briefly describe a leave one out (LOO) procedure for estimating the error of an SVM classifier on S . **2 pts**

(i) make n subset of S , each ~~one~~ exclude $\{\mathbf{x}_i, y_i\}$ in training set
and make it the validation set for \mathbf{x}_i

(ii) train ~~the~~ a SVM, for each \mathbf{x}_i , aim to minimize the training loss
(eg, misclassification error)

(iii) compute the error for \mathbf{x}_i , $\frac{1}{n} \sum_{i=1}^n I(\mathbf{c}(\mathbf{x}_i) \neq y_i)$

2. What is the LOO error? **1 pts**

$$\frac{1}{n} \sum_{i=1}^n I(\mathbf{c}(\mathbf{x}_i) \neq y_i) \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{c}(\mathbf{x}_i))^2$$

3. Suppose that we trained an SVM classifier on the **entire** dataset S , denote by sv the set of support vectors, $sv = \{\mathbf{x}_j | \alpha_j > 0\}$.

For the same value of C , prove that the LOO error is bounded by $\frac{|sv|}{n}$ i.e.

$$\text{LOO error} \leq \frac{|sv|}{n}$$

4 pts.

$$\text{Sign} \left(\sum_i \alpha_i y_i \mathbf{x}_i + b \right)$$

$$R^{CV} = \frac{1}{n} \sum_{i=1}^n I(\mathbf{c}(\mathbf{x}_i) \neq y_i)$$

$$\frac{1}{n} |TSV|$$

when $\mathbf{x}_i \notin sv$

we have

$$I(\alpha_i(\mathbf{x}_i) \neq y_i) = I(\alpha_i y_i = 0) = 0$$

when $\mathbf{x}_i \in sv$

$$I(\alpha_i(\mathbf{x}_i) \neq y_i) \leq 1$$

Question 3: Bagging and Boosting (20 pts.)

a) Answer precisely the following questions.

1. Are bagging and Boosting Bayesian approaches? Why?

No, they are heuristic approaches, 1 pt.

2. How is it possible to detect outliers with AdaBoost?

those who possess large weight are with the iteration of the computing progress 1 pt.

3. From the frequentist perspective, bagging is motivated by the tradeoff between two terms. Which?

Variance and bias 1 pt.

$$\left(\frac{\sum_{i=1}^B f_i(x)}{B} - E[f_b(x)] \right)^2 = \frac{B}{\sum_{i=1}^B} \left(f_b(x) - E[f_b(x)] \right)^2 = \frac{\sigma^2}{B}$$

4. AdaBoost has an alternative interpretation which is based on the minimization of a certain cost function. Which function?

1 pt.

$$E(e^{-yf(x)})$$

5. AdaBoost aims at selecting the best approximation to which ratio?

2 pt.

$$\text{prior ratio} \quad \epsilon_b = \frac{\sum_{i=1}^n w_i^{(b)} I(y_i \neq f_b(x_i))}{\sum_{i=1}^n w_i^{(b)}}$$

$$\hat{f}(x) = \frac{1}{2} \log \frac{P(y=1|x)}{P(y=-1|x)}$$

$$E_w(\alpha Y(x)/X)$$

$$E_w(\alpha Y(x))$$

$$e^{-\alpha Y(x)}$$

6. Why is the standard form of AdaBoost limited to binary classification?

$$G_b = \sum_{i=1}^n w_i \text{II}(Y_i t_i C_i X_i)$$

Since the form of $E_w(\alpha Y(x)/X)$ corresponds to the binary classifier 1 pt.

7. How could one parallelize bagging?

since bagging process different bootstrap dataset independently, so all single classifier could be trained at the same time parallelly, and the ensemble classifier just run the average. 1 pt.

8. Name a design property of the base classifiers of AdaBoost which impacts the overall predictive power.

weak classifier

the prediction accuracy should not be too high,
and should be higher than 0.5.

B. 1 pt.

9. Under which conditions does AdaBoost yield good results even when the base classifiers exhibit an individual performance that is only slightly better than that purely due to chance?

2 pt.

no too much outlier

Uncorrelated different classifier

B much larger than 1

and large dataset

- b) Consider *bagging* in the context of binary classification. Let the target function be $h(x)$, where $h : \mathbb{R}^d \rightarrow \{\pm 1\}$. Let us combine B individual classifiers $y_b(x)$, $b = 1 \dots B$ to obtain a committee model

$$y_{\text{COM}}(x) = \text{sign} \left[\frac{1}{B} \sum_{b=1}^B y_b(x) \right]. \quad (1)$$

1. Write down the pseudocode of bagging for binary classification, from the input (data and model) to the prediction output.

3 pts.

(1) bootstrap dataset to produce B bootstrap dataset for training
 original $\mathcal{D}_B \leftarrow \mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_B$.

(2) for $i = 1$ to B
 from classifiers $y_i(x)$ on \mathcal{D}_B separately
 individual

(3) combine the classifier and add equal weight to
 them, path to form the committee model
 $y_{\text{COM}}(x) = \text{sign} \left[\frac{1}{B} \sum_{b=1}^B y_b(x) \right]$

2. The error $\epsilon_b(x) = \exp\{-h(x)y_b(x)\}$ indicates the error of an individual model $y_b(x)$ for a single sample x in terms of the target function $h(x)$ and the output of the individual model $y_b(x)$. Write down E_{AV} , that is the average of the expected errors over the individual classifiers $y_b(x)$, and the expected error E_{COM} made by combined model $y_{\text{COM}}(x)$ as a function of the output of the committee model and of the target function.

$$E_{\text{AV}} = \frac{1}{B} \sum_{b=1}^B E[\exp(-h(x)y_b(x))]$$

$$E_{\text{COM}} = E[\exp \left\{ -h(x) \text{sign} \left[\frac{1}{B} \sum_{b=1}^B y_b(x) \right] \right\}] \quad 1 \text{ pts.}$$

3. Under which conditions is $E_{AV} < E_{COM}$?

$$\frac{\prod_{b=1}^B e^{y_b x}}{B} < e^{-\text{sign}\left(\sum_{b=1}^B y_b(x)\right)}$$

3 pts.

4. With the same exponential error, write down the error function for each iteration of AdaBoost with weighting coefficients α_b for the B base classifiers $y_b(x)$.

$$E_{\text{AdaBoost}} = e^{-Y(F^{(B)}(x) + \alpha_b y_b(x))}$$

1 pt.

5. In this scenario, within AdaBoost the minimization of this error function is performed with respect to two terms, which?

$$1) \quad y_b(x) = \begin{cases} 1 & p_w(y=1|x) > p_w(y=-1|x) \\ -1 & \text{otherwise} \end{cases}$$

$$2) \quad \alpha_b = \frac{1}{2} \ln \frac{1 - \epsilon_b}{\epsilon_b}$$

1 pt.

Question 4: Regression, Bias and Variance (20 pts.)

- a) Write down the linear regression model (component-wise and in vector notation) for input variable $x = (1, x_1, \dots, x_D)^T \in \mathbb{R}^{D+1}$ and output variable y . Formally introduce the model parameter(s).

$$y = x^T \beta + \epsilon$$

\downarrow \downarrow \downarrow
 $\epsilon \in \mathbb{R}^{D+1}$ $\beta \in \mathbb{R}^{D+1}$ $\epsilon \in \mathbb{R}$

From now on, assume that the input dataset consists of N samples given by the matrix $\mathbf{X} \in \mathbb{R}^{N \times (D+1)}$ (where the first column is $\mathbf{1}$), and the output, $\mathbf{Y} \in \mathbb{R}^N$. Write down the linear regression model for all observations \mathbf{Y} (in matrix notation).

2 pts.

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

- b) Write down the following cost functions (in a notation of your choice):

2 pts.

1. Ridge Regression (RR):

$$\sum_{i=1}^P (y_i - \sum_{j=1}^{D+1} x_{ij} \beta_j)^2 + \lambda \beta \sum_{j=1}^{D+1} \beta_j^2$$

$$(\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) + \lambda \beta^T \beta$$

2. Least Absolute Shrinkage and Selection Operator (LASSO).

Formulate as a constrained optimization problem:

$$\min \quad (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) \quad \min \sum_{i=1}^P |y_i - \sum_{j=1}^{D+1} x_{ij} \beta_j|^2$$

with $\|\beta\|_1 < \underline{\sum_{j=1}^{D+1} |\beta_j| \leq s}$

subject to

- c) Formulate the objective of **learning** (formula) in the regression setting introduced above with data (\mathbf{X}, \mathbf{Y}) i.i.d. from $P(\mathbf{X}, \mathbf{Y})$.

2 pts.

$$f = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- d) Assuming that the observations in \mathbf{Y} are affected by additive Gaussian noise ϵ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$. What do we know about the distribution of the *RSS*-estimator (i.e. $\hat{\beta}^{RSS}$)?

3 pts.

$$\hat{\beta}^{RSS} = E[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}] = \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E[\mathbf{X}\beta + \epsilon] \\ = \mathbf{a}^T \beta$$

$$V(\hat{\beta}^{RSS}) = V\{ \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \epsilon) \} \\ = \cancel{V(\mathbf{a}^T \beta)} \cdot V(\mathbf{a}^T \epsilon) + V(\mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon) \\ = V(\mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon) \quad \hat{\beta} \sim N(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}) \\ = \sigma^2 \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{a}$$

- e) Please briefly give another motivation for the *RSS* cost function.

1 pt.

according to Bayesian theory

assume the residual $\epsilon \sim N(0, \sigma^2)$

in order to maximize the likelihood

$$\max P(Y | X, \beta) \propto e^{-\frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta)}$$

$$\min (Y - X\beta)^T (Y - X\beta)$$

if give $P(\beta | \pi) \sim N(\mu, \Sigma)$

to maximize the posterior of β

$$P(\beta | Y, X, \pi) \propto e^{-\frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta) - \frac{1}{2} \beta^T \Sigma \beta}$$

$$\min \underbrace{(Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta}_{14}$$

Question 5: Unsupervised Learning (20 pts.)

1. In this section we study non-parametric density estimation of an arbitrary point x . We consider some small region \mathcal{R} containing x . In the class we have seen the following generic formula for density estimation:

$$p(x) = \frac{K}{nV},$$

where K denotes the number of data points falling inside the region \mathcal{R} and V shows the volume of the region. n is the number of data points in the sample set $\mathcal{S} = \{x_1, \dots, x_n\}$.

- (a) Consider the following Gaussian distribution to be used as a Parzen window function:

$$\phi(x - x_j) = \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{(x - x_j)^2}{2}\right). \quad (7)$$

What are K and V for this window function?

$$K = \sum_i \phi(\frac{x-x_i}{h})$$

$$V = 1$$

4 pts.

- (b) This particular choice of a window function leads to underfitting. Add a parameter to increase the model complexity.

too smooth, add h_n

such $p_{xy} = \frac{1}{n} \sum_{i=1}^n \phi\left(\frac{x-x_i}{h_n}\right)$

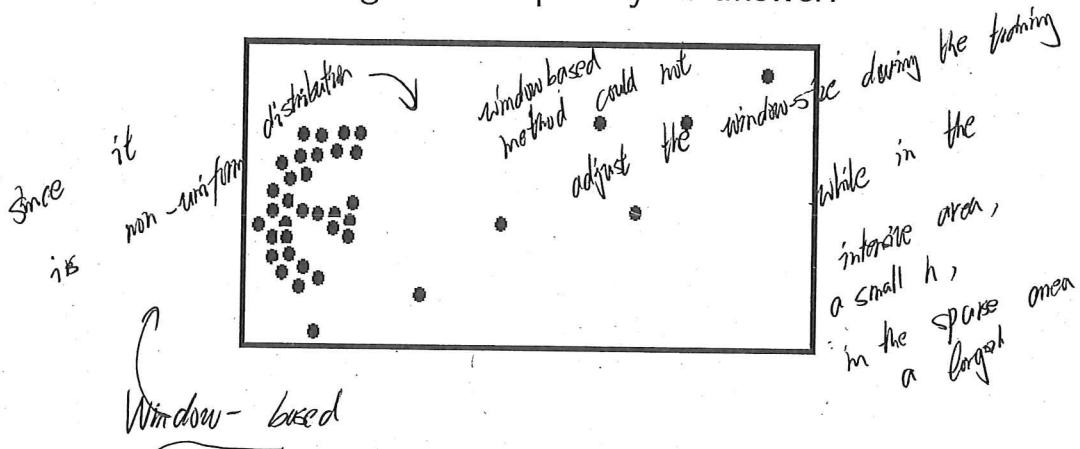
$$v_n = (h_n)^d$$

2 pts.

let h_n decrease, in this case.

the model will become more complex.

- (c) Consider the following sample set. Which of the density estimation methods would you choose? Window-based (Eq. (7)) or K -nearest neighbor? Explain your answer.



since it is obvious the data on the right is more sparse than the left one, thus possessing somehow different densities, then the K -NN would not take the density difference into account well as good as window-based is a way of hard assign so we choose the window-based, though

- (d) For a general Parzen window function prove that it provides a probability distribution.

for a parzen window

it meet the requirement

$$\text{that } \phi(x) > 0$$

$$\int \phi(x) dx = 1$$

4 pts.

2. We consider a mixture of K poisson distributions and perform the Expectation-Maximization (EM) algorithm to compute the unknown parameters. The log-likelihood function of n independent objects for mixture of K Poisson distribution is defined as:

$$P(x; \lambda) = \sum_{i=1}^n \log \sum_{c=1}^K \pi_c f(x_i; \lambda_c)$$

where π_c 's are the mixture weights and λ_c 's are the parameters of K Poisson distributions. $f(x_i; \lambda_c)$ is defined as:

$$f(x; \lambda_c) = \frac{\lambda_c^x e^{-\lambda_c}}{x!}$$

- (a) Introduce the latent indicator variables necessary for maximizing the log-likelihood function.

$M_{i,c}$ which ~~equals~~ equals 1 when x is assign to class c
otherwise 0

~~$$\sum_{i=1}^n \log \prod_{k=1}^K \pi_k^{M_{i,k}}$$~~

$$P(M_{i,c}=1) = \pi_c$$

$$P(x_i; \pi) = \prod_{k=1}^K \pi_k^{M_{i,k}}$$
2 pts.

- (b) Calculate the expectation of the latent variables. Provide a Bayesian interpretation for your answer.

$$E(M_{i,c}|x, \lambda) = P(\cancel{C}|x, \lambda) \cdot 1$$

$$= P(C|x, \lambda)$$

$$= \frac{P(c|x, \lambda) P(x|c, \lambda)}{P(x|\lambda)} = \gamma_{i,c}$$

$$\gamma_{i,c} = \frac{\pi_c P(m_{i,c}; \lambda_c)}{\sum_{k=1}^K \pi_k P(m_{i,k}; \lambda_k)}$$
5 pts.

(c) **Bonus question:** Assume the expectations of the latent variables are given. Calculate the unknown parameters λ_c 's. Write down the details of your calculations.

$$L(\theta, \pi_c | X) = \prod_{i=1}^n \sum_{c=1}^K \frac{\lambda_c^{x_{ic}} \pi_c}{x_i!} \quad x_{ic} \log \pi_c + \log \frac{\lambda_c^{x_{ic}}}{x_i!} - \lambda_c (x_i - 1)$$

$$L = \prod_{i=1}^n \sum_{c=1}^K \left[x_{ic} \log \pi_c + \log \frac{\lambda_c^{x_{ic}}}{x_i!} \right] - \lambda_c (x_i - 1)$$

5 pts.

$$\frac{\partial L}{\partial \lambda_c} = \prod_{i=1}^n x_{ic} \cdot \frac{\frac{x_{ic}-1}{x_i!} e^{-\lambda_c} + (-e^{-\lambda_c}) \lambda_c^{x_{ic}}}{\lambda_c^{x_{ic}} e^{-\lambda_c} / x_i!} = 0$$

$$\prod_{i=1}^n x_{ic} (x_i!) \left[\frac{x_{ic}}{x_i} - 1 \right] = 0$$

$$\sum_{i=1}^n \frac{x_{ic}}{x_i} = K \cdot \prod_{i=1}^n x_{ic}$$

$$\lambda_c = \frac{\sum x_{ic}}{\sum x_i}$$