

Where $\mathcal{L}(w_1, w_2, r)$ is the regularizer. We ask you to suggest a suitable regularizer.

(3)

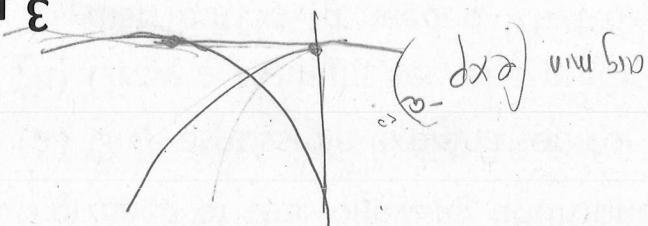
$$\text{argmin}_{w_1, w_2, r} \sum_i \exp(-y_i[r^2 - ((x_{i1} - w_1)^2 + (x_{i2} - w_2)^2)]) + \lambda \mathcal{L}(w_1, w_2, r)$$

new criterion will be

Figure 2 depicts the dataset we had, and h^* that we got using equation (2). To improve h^* we decided to add a regularizing term, the

tion (2).

3 pts.



$$\mathcal{L}(w_1, w_2, r) = r^2$$

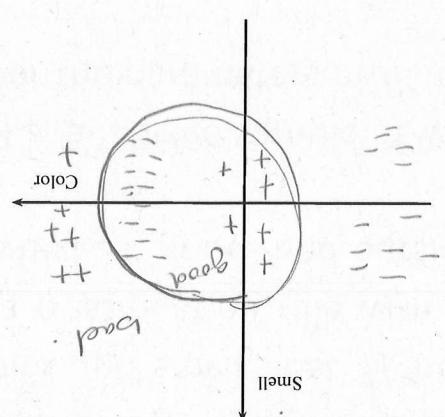
samples. You want to penalize the size of the circle.

2. Suggest a way to measure the empirical variance of the classifier h^* , given that we are out of budget for obtaining more watermelon samples.

8 pts.

Circle would never capture the structure of bimodal data

low training error



training error of ~ 0.4

Additional assumption: The sample size was large, and h^* had

4 pts.

Variance \uparrow 150

The size of the circle will decrease.

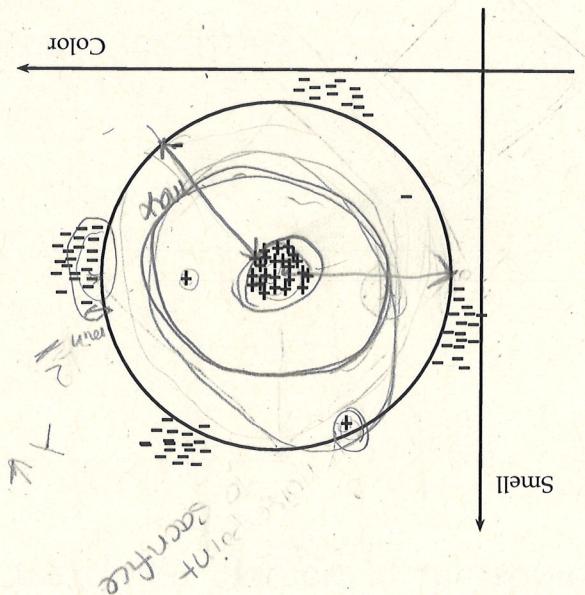
to some starting value $\lambda_0 > 0$.

what happens to the variance of h^* as we increase λ compared density regions visible in Figure 2, plus sparse outliers. Explain

5. Assume that the true distribution of watermelons consists of high

5 pts.

Figure 2: Watermelons dataset and h^*



$$\mathcal{D}(w_1, w_2, r) = R^2 \text{ in } l_2-\text{norm}$$

your answer.

(b) Write down the mathematical term of the regularizer. Explain

4. (a) Draw on Figure 2 the regularized solution you envision.

$$\frac{\partial}{\partial w} = y_i - w^T \phi(x_i) - b = 0 \quad \text{as minimization problem}$$

- the features in the kernel space and rewrite the problem as a
- (a) We replace the inputs x_i in Equation (4) with the vectors of

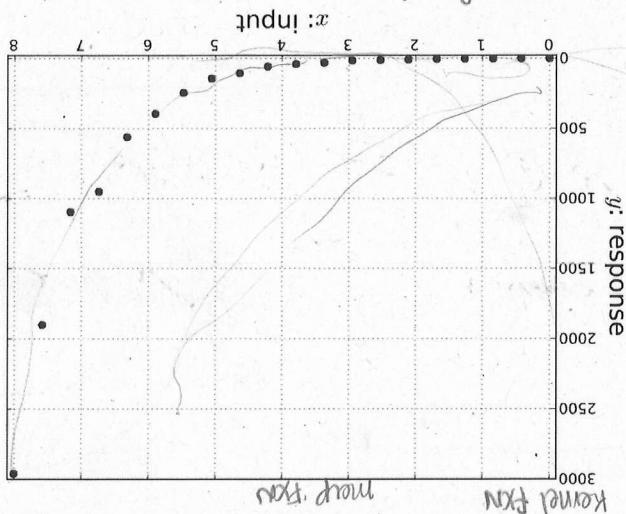
$$\min_w \sum_i (y_i - w^T \phi(x_i) - b)^2 + \frac{\lambda}{2} \|w\|^2 \quad (4)$$

Recall the formulation of ridge regression as an optimization problem:

- linear regression solution, for datasets such as the one depicted introducing a feature transform $\Phi(x)$. This should allow a non-
2. You will now derive a kernelized version of ridge regression by

$K(x_i, x_j) = \dots$

kernel method used for mapping from to make data linearly separable



$$k(x) = \log(x_i) \cdot \log(x_j)$$

$$\phi(x) = \log(x)$$

space. defn: $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ $\phi(x) \rightarrow \text{Kernel space}$

- such that you can use a linear regression method in the kernel space variable y . Your task is to find a kernel function $K(x_i, x_j)$, and response variables y_i , the goal is to find a functional relation between them, often expressed with a weight vector w and bias b .
- Below is a dataset with one dimensional input variables x , and response variables y_i , the goal is to find a functional relation between them, often expressed with a weight vector w and bias b .

Recall the regression setting: Given input vectors x_i , and output (response) variables y_i , the goal is to find a functional relation between them, often expressed with a weight vector w and bias b .

Question 4: Kernelized Ridge Regression (20 pts.)

8 pts.

$$L(\alpha) = -\frac{1}{2} \sum_i \alpha_i^2 + \sum_i \alpha_i y_i - \frac{1}{2} \sum_i \alpha_i \phi_i^\top \phi_i \quad \text{min } L(\alpha)$$

$$\|W\|^2 = \frac{1}{2} \cdot \frac{1}{2} \|\alpha \phi\|^2 = \sum_i \alpha_i \phi_i^\top \phi_i$$

you want to get the objective function of alpha

$$\frac{\partial L}{\partial \alpha} = 2\alpha - \alpha_i = 0 \quad \Leftrightarrow \quad \alpha_i = \frac{1}{2}$$

$$\frac{\partial L}{\partial W} = -\sum_i \alpha_i \phi_i + \lambda W = 0 \quad \Leftrightarrow \quad W = \frac{1}{\lambda} \sum_i \alpha_i \phi_i$$

this is why you derive on α .

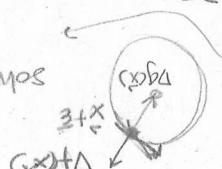
$$L(\alpha, W, b, \phi) = \frac{1}{2} \alpha^\top \alpha + \frac{1}{2} \|W\|^2 + \sum_i \alpha_i (y_i - \phi^\top W - b) - \frac{1}{2} \lambda$$

(c) Derive the dual optimization problem.

$$\frac{\partial L}{\partial \alpha} = g(\alpha) = 0$$

$$\frac{\partial L}{\partial x} = 0$$

solve the lagrangian



$$s.t. \quad g(\alpha) = 0 \quad L(\alpha, x) = f(x) + g(x)$$

$$\text{max } f(x)$$

objective of lagrangian is to minimize the problem

using α as the dual variable.

(b) Write down the Lagrangian of this new optimization problem

2 pts.

$$(6) \quad \dots = \dots \quad \text{s.t.}$$

$$(5) \quad \min_w \sum_i \xi_i^2 + \frac{1}{2} \|w\|^2$$

penalty

ables ξ_i . Write down the equality constraint in Equation (6). constrained optimization problem by introducing the new variables

4 pts.

$$= \sum_i \alpha_i * \phi_i^T \phi_k = \sum_i \alpha_i * K(x_i, x_k)$$

$$\text{Plug in } w = \sum_i \alpha_i \phi_i$$

$$y_k = \sum_i \phi_i^T \phi_k$$

point x_k , write down the equation to compute y_k .

4. Given the optimal solution of the dual problem α^* and a new

3 pts.

$$K(x_i, x_j) = \phi_i^T \phi_j$$

3. Express the dual problem in terms of the kernel function $K(x_i, x_j)$.

2 pts.

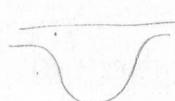
$$\phi\left(\frac{(x-x_j)^2}{h^2}\right) = \frac{1}{\sqrt{2\pi}h} \exp\left(-\frac{(x-x_j)^2}{2h^2}\right)$$

Add variance parameter $= h$, controls smoothness

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

4 pts.

$V = 1 \rightarrow$ because a probability distribution



$$K = \sum_{j=1}^n \phi(x-x_j)$$

What are K and V for this window function?

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

Parzen window function:

- (a) Consider the following Gaussian distribution to be used as a

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

$$d(x) = \frac{nV}{K}$$

density estimation:

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

Question 5: Unsupervised Learning (20 pts.)

4 pts.

$$\frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\frac{1}{n} \sum_{j=1}^n \phi(x - x_j)$$

$$\frac{1}{(2\pi)^{\frac{n}{2}}} |\mathcal{E}|^{\frac{n}{2}}$$

$$P(x) = \frac{k}{nV} = \frac{1}{n} \sum_{j=1}^n \phi(x - x_j) \Leftrightarrow \int p(x) dx = 1$$

due to a probability distribution $\Phi(u) = \int \phi(u) du$

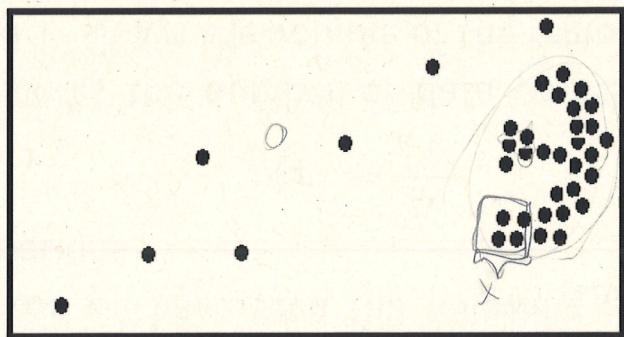
must satisfy two conditions $\Phi(u) \geq 0$

probability distribution.

(d) For a general Parzen window function prove that it provides a

3 pts.

choose K -nearest neighbor. \Rightarrow local-based



or K -nearest neighbor? Explain your answer.

(c) Consider the following sample set. Which of the density estimation methods would you choose? Window-based (Eq. (7))

5 pts.

$$\sum_{j=1}^k P(x_n | z_n, \alpha_j) = \prod_{j=1}^k P(x_n | \alpha_j)$$

$$E[z_{nk}] = P(z_{nk}=1 | x_n) * 1 = \prod_{j=1}^k P(z_{nj}=1) = P(x_n | \alpha_k)$$

Joint dist're, button

$$M\text{-step: } \theta_{new} = \arg \max_{\theta} Q(\theta, \theta_{old})$$

$$P(z|x, \theta_{old}) = E_z(z|\dots)$$

E-step: estimate posterior distribution of data

Bayesian interpretation for your answer.

(b) Calculate the expectation of the latent variables. Provide a

2 pts.

Conditional expectation
Expectation

$$z_{nk} = \begin{cases} 0 & : \text{aw} \\ 1 & : \text{component R is responsible for generated } x_n \end{cases}$$

$\vec{z} \in \mathbb{R}^{N \times K}$ each row is a vector for a point

$$\log P(x|\vec{z}, \vec{\alpha}) = \sum_{n=1}^N \log \prod_{k=1}^K f(x_n | \alpha_k)$$

using the log-likelihood function.

(a) Introduce the latent indicator variables necessary for maximization.

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

of K Poisson distributions: $f(x_i; \lambda_i)$ is defined as:

where π_c 's are the mixture weights and λ_c 's are the parameters

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

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:

5 pts.

Pg 445, 430

$$\begin{aligned}
 & \sum_{k=1}^n \sum_{i=1}^{x_k} \log p(x_i | z_i, \theta) = \\
 & = \left[\prod_{i=1}^n \prod_{k=1}^{x_i} p(x_i | z_i, \theta) \right] \\
 & \text{Objective fun } Q = E \left[\log p(x, z | \theta) \right]
 \end{aligned}$$

down the details of your calculations.

ables are given. Calculate the unknown parameters θ 's. Write

(c) **Bonus question:** Assume the expectations of the latent vari-

conditional expectations

Q/A session before exam.

SVM \rightarrow assume vertical space behind ~~data~~

CV scheme is really important

choose solution for insights and do it better \rightarrow do not remove

the only goal.

No, fixed "right solution", as in real life. \hookrightarrow The accuracy is

