

Answer ALL SIX questions. You may use results from the notes without reproving them, however, add a citation.

Notation: Let $[m] := \{1, \dots, m\}$. We also overload notation so that

$$[\text{pred}] := \begin{cases} 1 & \text{pred is true} \\ 0 & \text{pred is false} \end{cases}.$$

Marks for each part of each question are indicated in square brackets.

1. We consider problem of empirical risk (error) minimisation for multivariate simple polynomial functions. For all $\alpha \in \mathbb{R}^n$ define $h_\alpha : (0, \infty)^n \rightarrow (0, \infty)$ as

$$h_\alpha(\mathbf{x}) = x_1^{\alpha_1} x_2^{\alpha_2} \times \dots \times x_n^{\alpha_n},$$

the hypothesis space,

$$\mathcal{H}_{\text{exp}} := \{h_\alpha : \alpha \in \mathbb{R}^n\}$$

and the error function $\ell(y, \hat{y}) := \log(\frac{y}{\hat{y}})^2$. Given a dataset $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \subset (0, \infty)^n \times (0, \infty)$ design an efficient (polynomial-time) algorithm to perform empirical risk minimisation with respect to \mathcal{H}_{exp} and ℓ . Argue that your algorithm is correct.

[10 marks]

2. a. Suppose $K : [0, \mathbb{R}) \times [0, \mathbb{R}) \rightarrow \mathbb{R}$ is a kernel.

- i. Is $K(x^3, t^3)$ a kernel? Explain why or why not.
- ii. Is $K(\sqrt{x}, \sqrt{t})$ a kernel? Explain why or why not. Note for the purpose of this question $\sqrt{\cdot}$ is single-valued and returns the positive root.

[5 marks]

b. A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ is *generic* iff $\mathbf{x}_i = \mathbf{x}_j \Rightarrow y_i = y_j$. Consider the following kernel function $K_a(\mathbf{x}, \mathbf{t}) := \prod_{i=1}^n (1 + (x_i t_i) + (1 - x_i)(1 - t_i))$.

- i. Argue that $K_a : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a kernel.
- ii. Given a generic training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \{0, 1\}^n \times \{0, 1\}$ does there necessarily exist a vector $\alpha \in \mathbb{R}^m$ such that

$$\sum_{j=1}^m \left(\sum_{i=1}^m \alpha_i K_a(\mathbf{x}_i, \mathbf{x}_j) - y_j \right)^2 = 0?$$

Provide an argument to justify your answer.

[5 marks]

3. Linear support vector machines (SVMS)

Assume that the set $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \subset \mathbb{R}^2 \times \{-1, 1\}$ of binary examples is strictly linearly separable by a line going through the origin, that is, there exists a vector $\mathbf{w} \in \mathbb{R}^2$ such that the linear function $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, $\mathbf{x} \in \mathbb{R}^2$ has the property that $y_i f(\mathbf{x}_i) > 0$ for every $i = 1, \dots, m$. Consider the optimisation problem (linearly separable SVM):

$$\mathbf{P1} : \quad \text{minimise } \left\{ \frac{1}{2} \mathbf{w}^\top \mathbf{w} : y_i \mathbf{w}^\top \mathbf{x}_i \geq 1, i = 1, \dots, m \right\}.$$

- a. Argue that the above problem has a unique solution. Describe the geometric meaning of this solution.

[5 marks]

- b. Show that the vector \mathbf{w} solving problem **P1** has the form $\mathbf{w} = \sum_{i=1}^m c_i y_i \mathbf{x}_i$ where c_1, \dots, c_m are some nonnegative coefficients. [HINT: use the method of Lagrange multipliers]

[5 marks]

- c. Show that the coefficients c_1, \dots, c_m in the above formula solve the optimization problem

$$\mathbf{P2} : \quad \max \left\{ -\frac{1}{2} \sum_{i,j=1}^m c_i c_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^m c_i : c_j \geq 0, j = 1, \dots, m \right\}.$$

Finally, if $(\hat{c}_1, \dots, \hat{c}_m)$ is the solution to this problem and $\hat{\mathbf{w}}$ is the solution to problem **P1**, argue that $\hat{\mathbf{w}}^\top \hat{\mathbf{w}} = \sum_{i=1}^m \hat{c}_i$.

[5 marks]

[Question 3 cont. over page]

[Question 3 cont.]

- d. Now drop the linear separability assumption. Consider the following data dependent representation: $\mathbf{z}_i = (\mathbf{x}_i, \rho \mathbf{e}_i) \in \mathbb{R}^{d+m}$, where ρ is a positive parameter and \mathbf{e}_i is the m dimensional vector all of whose components are equal to zero except for the i -th component which is equal to 1. Argue that the binary data $\{(\mathbf{z}_i, y_i)\}_{i=1}^m$ are linearly separable. Furthermore show that the corresponding problem,

$$\mathbf{P3} : \min_{\mathbf{v} \in \mathbb{R}^{2+m}} \left\{ \frac{1}{2} \mathbf{v}^\top \mathbf{v} : y_i \mathbf{v}^\top \mathbf{z}_i \geq 1, i = 1, \dots, m \right\}$$

is equivalent to the problem

$$\mathbf{P4} : \min_{\mathbf{w} \in \mathbb{R}^2, \xi \in \mathbb{R}^m} \left\{ \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \frac{1}{2\rho^2} \sum_{i=1}^m \xi_i^2 : y_i \mathbf{w}^\top \mathbf{x}_i \geq 1 - \xi_i, i = 1, \dots, m \right\}.$$

Then provide an interpretation of the latter problem as regularisation problem, indicating which is the loss function used. (2-5 sentences recommended.)

[5 marks]

4. Adaboost.

Recall the following notation,

- α_t : the weight on weak learner t where $\alpha_t \in \mathbb{R}$.
- $D_t(i)$: the weight on example i at time t where $\sum_{i=1}^m D_t(i) = 1$
- ϵ_t : “weighted error of weak learner $h_t(\cdot)$ at time t ”

$$\epsilon_t := \sum_{i=1}^m D_t(i) [h_t(\mathbf{x}_i) \neq y_i]$$

- training error at time T ,

$$\frac{1}{m} \sum_{i=1}^m [H(\mathbf{x}_i) \neq y_i]$$

- Briefly discuss the role of the “ α ” weights and “ D ” weights in Adaboost (2-5 sentences recommended). Your discussion should remark on the significance of relatively high or low weight values.

[5 marks]

- Let $\epsilon^* = \max_{t \in [T]} \epsilon_t$, use ϵ^* to bound the training error at time T .

[5 marks]

- Look-ahead weighted error.*

$$\hat{\epsilon}_t := \sum_{i=1}^m D_{t+1}(i) [h_t(\mathbf{x}_i) \neq y_i],$$

What, if anything, can we infer about the $\hat{\epsilon}_t$ when $t = T$. Explain your reasoning.

[10 marks]

5. The problem of designing efficient algorithms that obtain (expected) non-trivial regret bounds for linear classification with respect to the “0-1” loss (misclassifications) seems to be a difficult problem. In this question we consider simpler problems.

We have the following definitions the 0-1 loss is $\ell_{01}(y, \hat{y}) := [y \neq \hat{y}]$, the hypothesis class of 2-norm bounded linear *classifiers* over a set $\mathcal{X} \subset \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| = 1\}$,

$$\hat{\mathcal{H}}_{\mathcal{X},U} = \{h_{\mathbf{u}} : \mathbf{u} \in \mathbb{R}^n, \|\mathbf{u}\| \leq U, \forall \mathbf{x} \in \mathcal{X} : |\mathbf{u} \cdot \mathbf{x}| \geq 1\}$$

where $h_{\mathbf{u}} : \mathcal{X} \rightarrow \mathbb{R}$ and $h_{\mathbf{u}}(\mathbf{x}) := \text{sign}(\mathbf{u} \cdot \mathbf{x})$, and the corresponding hypothesis class of linear *interpolants*

$$\bar{\mathcal{H}}_{\mathcal{X},U} = \{h_{\mathbf{u}} : \mathbf{u} \in \mathbb{R}^n, \|\mathbf{u}\| \leq U, \forall \mathbf{x} \in \mathcal{X} : |\mathbf{u} \cdot \mathbf{x}| = 1\},$$

intuitively the functions in $\hat{\mathcal{H}}$ classify the points in \mathcal{X} with a margin ≥ 1 and the functions in $\bar{\mathcal{H}}$ interpolate the points in \mathcal{X} to be exactly 1 or -1. Also observe that $\bar{\mathcal{H}}_{\mathcal{X},U} \subseteq \hat{\mathcal{H}}_{\mathcal{X},U}$.

- a.
 - i. Discuss why giving regret bounds for predicting linear classifiers with respect to the 0-1 loss may be more difficult than with the hinge loss. Please limit your discussion to 1-2 sentences.
 - ii. Discuss why from a “user’s perspective” it may be more useful to have a 0-1 loss than a hinge loss bound. Please limit your discussion to 1-3 sentences.

[5 marks]

[Question 5 cont. on next page]

[Question 5 cont.]

b. *Simplifying by restricting the hypothesis class.*

Protocol:

Nature selects $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \in \mathcal{X}$ and $y_1, y_2, \dots, y_T \in \{-1, 1\}$.

For $t = 1$ To T Do

Receive pattern	$\mathbf{x}_t \in \mathcal{X}$
Predict	$\hat{y}_t \in \{-1, 1\}$
Receive label	$y_t \in \{-1, 1\}$

Design a polynomial-time randomised algorithm with a good upper bound on the expected regret,

$$\mathbb{E} \left[\sum_{t=1}^T \ell_{01}(y_t, \hat{y}_t) \right] - \min_{h \in \mathcal{H}_{\mathcal{X}, U}} \sum_{t=1}^T \ell_{01}(y_t, h(\mathbf{x}_t))$$

give your algorithm design and argument for the upper bound.

[5 marks]

[Question 5 cont. over page]

[Question 5 cont.]

- c. *Simplifying by loosening the definition of regret.* In this part we generalise regret to allow a leading term $c(U)$ in front of the loss of the best hypothesis this leading term can only depend on U .

Protocol:

Nature selects $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \in \mathcal{X}$ and $y_1, y_2, \dots, y_T \in \{-1, 1\}$.

For $t = 1$ To T Do

Receive pattern	$\mathbf{x}_t \in \mathcal{X}$
Predict	$\hat{y}_t \in \{-1, 1\}$
Receive label	$y_t \in \{-1, 1\}$

Design a polynomial-time algorithm with a good upper bound on the $c(U)$ -regret,

$$\sum_{t=1}^T \ell_{01}(y_t, \hat{y}_t) - c(U) \min_{h \in \hat{\mathcal{H}}_{\mathcal{X}, U}} \sum_{t=1}^T \ell_{01}(y_t, h(\mathbf{x}_t)).$$

Observe that the introduction of $c(U)$ generalises the notion of regret where $c(U) = 1$ recovers the usual notion of regret. Thus the aim is that $c(U)$ is small while still guaranteeing that,

$$\frac{\sum_{t=1}^T \ell_{01}(y_t, \hat{y}_t)}{T} \leq c(U) \min_{h \in \hat{\mathcal{H}}_{\mathcal{X}, U}} \frac{\sum_{t=1}^T \ell_{01}(y_t, h(\mathbf{x}_t))}{T} \text{ as } T \rightarrow \infty.$$

Give your algorithm design and argument for the upper bound.

[10 marks]

6. a. Given an example of a hypothesis class \mathcal{H} with $|\mathcal{H}| = \infty$ and $\text{vcdim}(\mathcal{H}) = 1$. Argue that it has $\text{vcdim}(\mathcal{H}) = 1$.

[5 marks]

- b. Suppose we have hypothesis class \mathcal{H} with $|\mathcal{H}| = \infty$ and $\text{vcdim}(\mathcal{H}) = 1$. What does this imply, if anything, about the possible existence of an online algorithm with a mistake bound for this hypothesis class. Explain your reasoning.

[5 marks]

- c. Suppose we have hypothesis class \mathcal{H} and an algorithm with a mistake bound of B for \mathcal{H} . What, if anything, does this imply about the pac-learnability of the hypothesis class. Explain your reasoning.

[10 marks]