

Practice problems to prepare for Exam 1

COMS 4771 Fall 2016

Problem 1 (Linear classifiers). Which of the following classifiers are *linear classifiers* in \mathbb{R}^d ?

- The function $f: \mathbb{R}^d \rightarrow \{0, 1\}$ given by $f(\mathbf{x}) = \mathbb{1}\left\{\frac{1}{1+\exp(-\langle \mathbf{w}, \mathbf{x} \rangle)} - \frac{1}{2} > 0\right\}$ for some vector $\mathbf{w} \in \mathbb{R}^d$.
- The 1-NN classifier based on the following training data ($d = 2$):

	Feature 1	Feature 2	Label
Example 1	-1	-1	-1
Example 2	+1	+1	+1
Example 3	-1	+3	-1
Example 4	+1	+5	+1

- A classifier trained using Kernelized Online Perceptron with the kernel $K(\mathbf{x}, \tilde{\mathbf{x}}) := \langle \mathbf{x}, \mathbf{A}\tilde{\mathbf{x}} \rangle$, where \mathbf{A} is a $d \times d$ symmetric positive definite matrix.

Problem 2 (Maximum likelihood estimation).

- (a) Consider the model \mathcal{P} of distributions over the positive integers \mathbb{N} with probability mass functions given by

$$P_\theta(x) = (1 - \theta)^{x-1}\theta, \quad x \in \mathbb{N}$$

for parameter $\theta \in (0, 1)$ called the *success parameter*.

Derive a formula for the maximum likelihood estimator for the success parameter given data $\{x_i\}_{i=1}^n$ (treated as an iid sample).

- (b) Consider the model \mathcal{P} of probability distributions over the non-negative reals \mathbb{R}_+ with probability densities on the given by

$$p_\lambda(x) = \lambda e^{-\lambda x}, \quad x \in \mathbb{R}_+$$

for parameter $\lambda > 0$ called the *rate parameter*.

Derive a formula for the maximum likelihood estimator for the rate parameter given data $\{x_i\}_{i=1}^n$ (treated as an iid sample).

Problem 3 (Kernels). Let $K_1: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and $K_2: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be kernel functions.

- (a) Define $K_3: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$K_3(\mathbf{x}, \mathbf{x}') := K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}').$$

Is K_3 always a kernel function (whenever K_1 and K_2 are)? Give a brief justification for your answer.

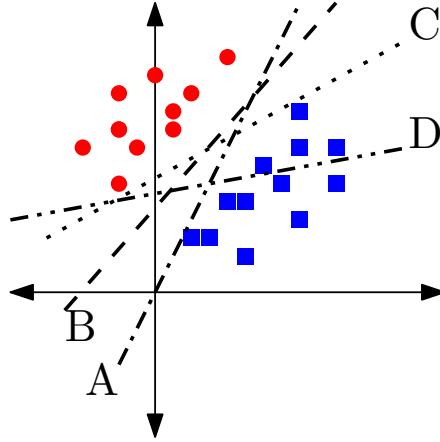
- (b) Define $K_4: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$K_4(\mathbf{x}, \mathbf{x}') := -K_1(\mathbf{x}, \mathbf{x}') - K_2(\mathbf{x}, \mathbf{x}').$$

Is K_4 always a kernel function (whenever K_1 and K_2 are)? Give a brief justification for your answer.

- (c) Let S be a data set of n labeled examples from $\mathcal{X} \times \{-1, 1\}$. Furthermore, let $\phi: \mathcal{X} \rightarrow \{z \in \mathbb{R}^D : \|z\|_2 = 1\}$ be the feature mapping corresponding to the kernel $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Suppose the shortest weight vector $\mathbf{w} \in \mathbb{R}^D$ such that $y\langle \mathbf{w}, \phi(\mathbf{x}) \rangle \geq 1$ for all $(\mathbf{x}, y) \in S$ has length $\|\mathbf{w}\|_2 = 10$. At most how many times does Kernelized Online Perceptron compute the value of the kernel function K when run on input S ? Give as small of an upper bound as possible, and briefly justify your answer.

Problem 4 (Linear classifiers, again). The figure below depicts the decision boundaries of four linear classifiers (labeled A, B, C, D), and the locations of some labeled training data (positive points are squares, negative points are circles).



Consider the following algorithms for learning linear classifiers which could have produced the depicted classifiers given the depicted training data:

1. (Batch) Perceptron
2. Online Perceptron (making one pass over the training data)
3. ERM for homogeneous linear classifiers
4. An algorithm that exactly solves the SVM problem

What is the most likely correspondence between these algorithms and the depicted linear classifiers? Give a brief explanation for your matching.

Problem 5 (True/false). Determine if each of the following statements is true or false. (For your own edification, you should also come up with a brief justification for your answer.)

- (a) For any classifier $f: \mathcal{X} \rightarrow \{0, 1\}$, if S is an iid sample from a distribution P over $\mathcal{X} \times \{0, 1\}$, then $\text{err}_S(f)$ is an unbiased estimator of $P(f(X) \neq Y)$, where $(X, Y) \sim P$.
- (b) It is possible for a leaf in a decision tree (during greedy training) to have zero “classification error” but non-zero “Gini index”.
- (c) Suppose an iid sample from a distribution P over $\mathcal{X} \times \{0, 1\}$ is randomly partitioned into three sets, S , V , and T . Further, suppose a feature map $\phi: \mathcal{X} \rightarrow \mathbb{R}^D$ is determined based on $S \cup V$, and a classifier $\hat{f}: \mathbb{R}^D \rightarrow \{0, 1\}$ is trained using $\{(\phi(\mathbf{x}), y) : (\mathbf{x}, y) \in S\}$. Then $\text{err}_V(\hat{f} \circ \phi)$ is an unbiased estimator of the true error rate of f .

- (d) Consider the setting from (c) through the penultimate sentence. Then $\text{err}_T(\hat{f} \circ \phi)$ is an unbiased estimator of the true error rate of f .
- (e) A classifier based on a generative model with Gaussian class conditional distributions is a linear classifier if the class conditional distributions are *product distributions*.
- (f) The linear classifier with the smallest training error rate on a set of examples S can be obtained as the solution to a linear program.
- (g) If Perceptron is run on a linearly separable data set S , and \mathbf{w}_\star is the shortest vector satisfying $y\langle \mathbf{w}_\star, \mathbf{x} \rangle \geq 1$ for all $(\mathbf{x}, y) \in S$, then the angle between weight vector $\hat{\mathbf{w}}_t$ maintained by Perceptron and \mathbf{w}_\star decreases after every update.
- (h) Every kernel function $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies $K(\mathbf{x}, \mathbf{x}') \geq 0$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$.
- (i) Suppose $(\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_n)$ is the solution for the SVM dual problem when the (linearly separable) training data is $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$. Let $S' := \{(\mathbf{x}_i, y_i)\}_{i=1}^{n+m}$ be a superset of S , and let $(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_{n+m})$ be the solution for the SVM dual problem when the training data is S' . For each $i = 1, 2, \dots, n$, whenever $\hat{\alpha}_i = 0$, then $\hat{\beta}_i = 0$. That is, if (\mathbf{x}_i, y_i) is not a support vector for the SVM classifier with data set S , then it is not a support vector for the SVM classifier with data set S' .
- (j) The soft-margin SVM optimization problem has no solution whenever the data set S contains a pair of examples (\mathbf{x}, y) and (\mathbf{x}, y') where $y \neq y'$.