

CSE517a Final

Good luck!

May 2015

NAME:	
Student ID:	
Email:	

General ML	
Kernls	
Neural Networks	
GPs	
Unsupervised Learning	
TOTAL	

1 [??] General Machine Learning

Either circle **T** or **F**. Questions declared as **True** require no explanation (worth 1 point). Questions declared as **False** require a **one sentence** explanation (worth 2 points).

T/F The fundamental difference between Bayesian and Frequentist Statistics boils down to the definition of good pop music. (For example, “Bayesians” love Meghan Trainor’s song “All about the Bayes” and just cannot hear enough of it, whereas “Frequentists” believe strongly that this song has been played way too frequently on the media in recent months. Frequentists are also concerned that certain passages may be inappropriate, in particular when the singer refers to *shaking her posterior*, which should be changed to *likelihood*.)

T/F The difference between Parametric and Non-parametric algorithms is that non-parametric algorithms have no hyper-parameters to tune.

T/F Platt scaling is used to scale features to be between 0 and 1.

T/F One advantage of 1 *vs.* 1 multi-class classification over 1 *vs.* *all* is that the problems are more likely to be class balanced.

T/F After convergence, the K-means cluster centers are always original data points.

T/F SVD decomposes a matrix \mathbf{X} (with data as column vectors) into three terms $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$. The matrix \mathbf{S} is diagonal and non-negative, the matrix \mathbf{U} is the projection matrix of PCA the matrix \mathbf{V} is the whitened data.

T/F The i^{th} largest eigenvalue of the covariance corresponds to the amount of variance captured in the i^{th} principal component.

T/F The Support Vector Machine with the RBF kernel is non-parametric.

T/F K-means is a non-parametric algorithm

T/F Kernel SVM (in the dual) with a linear kernel is always slower than training the linear SVM directly in the primal (without a kernel matrix), because the kernel matrix has size $n \times n$.

T/F Any real *symmetric* matrix is a well defined kernel.

T/F A sign for a high variance scenario is that the training error is higher than the testing error.

T/F An effective way to fight a **high bias** scenario is to substantially increase the amount of training data.

T/F Performing leave-one-out cross validation with the squared loss requires you to train n classifiers. Each time you leave out one point, train the classifier on the remaining $n - 1$ points and compute the error of the resulting classifier on the left out point. The error is the average error across all n classifiers.

2 [20] Kernels

1. (5) Assume Ω is a set of all words in the English dictionary, $|\Omega| = d$. We define a text document as the set of all the words that are in the document, $S \subseteq \Omega$. Proof that the following kernel over such documents is well-defined (*i.e.* positive semi-definite)

$$k(S_1, S_2) = \exp \left(\frac{|S_1 \cap S_2|}{|S_1| |S_2|} \right).$$

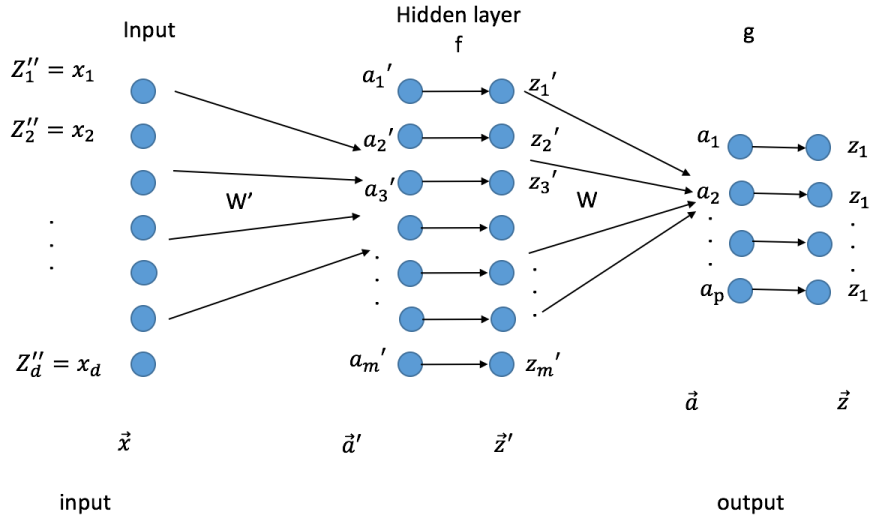
(Hint: It helps to write $\Omega = \{\omega_1, \dots, \omega_n\}$ and represent a set S as a binary vector $x \in \{0, 1\}^d$, where the i^{th} dimension $x_i = 1$ if and only if $\omega_i \in S$.)

2. (10) Assume you are given a data set $D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$ with $\vec{x}_i \in \mathcal{R}^d$ and $y_i \in \{+1, -1\}$. The Perceptron algorithm learns a separating hyper-plane in the following way:
 - (a) Initialize $\mathbf{w} = \vec{0}$
 - (b) Pick (\vec{x}_i, y_i) randomly from D .
 - (c) If $y_i \mathbf{w}^\top \vec{x}_i \leq 0$ then make the update $\mathbf{w} \leftarrow \mathbf{w} + y_i \vec{x}_i$, otherwise do nothing.
 - (d) Goto step (b) and repeat until convergence.

Show that this algorithm can be kernelized. First show that \mathbf{w} can be written as a linear combination of the inputs. State the kernelized classifier $h()$ and the kernelized version of the learning algorithm.

3. (5) Joe wants to combine the power of kernel machines with that of neural networks. He trains a neural network with one hidden layer on his training data. Then he defines the mapping $f(\vec{x}) = \vec{z}'$ to be the mapping of his input to the first hidden layer (after the transition function has been applied). (See the neural network figure on the next page for an illustration.) He then defines his kernel matrix $k(\vec{x}_1, \vec{x}_2) = f(\vec{x}_1)^\top f(\vec{x}_2)$ to train an SVM. Is this kernel function well defined?

3 [18] Neural Networks



- (3) Write z and z' in terms of x, a', a, W', W and the transition functions f and g .

- (5) Given $\delta_j = \frac{\partial L}{\partial a_j}$, derive $\delta'_j = \frac{\partial L}{\partial a'_j}$.

4 [11] Gaussian Processes

1. (5) What are the assumptions of Gaussian Process Regression (with mean $\vec{\mu}_x$ and variance \mathbf{K}_x) about the distribution of the labels y_1, \dots, y_n and a test label y_t ?

2. (5) The mean prediction of GPR is identical to kernel regression. So why can't we easily use kernel regression for hyper-parameter optimization with the Upper Confidence Bound (UCB) exploration / exploitation strategy?

3. (5) You are using Bayesian Global Optimization to optimize the hyper-parameters of an SVM (σ and C) with Upper Confidence Bound (UCB). After a small number of training points have been sampled your model predicts that the setting σ_m, C_m yields the highest predicted accuracy on the validation set. Explain a scenario where your optimization algorithm would pick a different set of hyper-parameters σ', C' to explore with lower predicted accuracy.

5 [19] Unsupervised Learning

1. (4) Batman wants to use k -means to cluster a data set of n data points. He doesn't know how many clusters k the data set has so he tries every possible setting $k = 1, \dots, n$ and computes the squared reconstruction loss on his training data for each one. Then his algorithm selects the clustering with the lowest reconstruction error and outputs it as the final answer. After the algorithm is done he is amazed to note that the best reconstruction error is in fact *zero*. He sees this as a clear indication that his data has strong cluster structure, which he has now uncovered. Do you share his opinion? What value of k did the algorithm select?

2. (5) In the K-means optimization algorithm, we use the following update rule for each cluster where $\vec{\mu}_i$ is the center of the i -th cluster.

$$\gamma_{ij} = \begin{cases} 1 & : j = \min_j (\vec{x}_i - \vec{\mu}_j)^2 \\ 0 & : \text{else} \end{cases}$$

$$\vec{\mu}_j = \frac{1}{\sum_i \gamma_{ij}} \sum_i \gamma_{ij} \vec{x}_i$$

Prove that this minimizes the squared Euclidean distance between each point \vec{x}_i and its cluster center $\vec{\mu}_j$.

