# Machine Learning 2017
# Final Exam

## 13 December 2017

Name: ...............................................................

**Question 1:** $\boxed{\text{2 points}}$ Consider a regression task in one dimension $x_1$, and let $\mathcal{D} = \{(0, 1), (1, 2)\}$ be a dataset of input-output pairs on the form $(x_1, y)$. If we perform linear regression, what is the optimal weight vector $\mathbf{w}_{\text{lin}}$? Recall that each input $(x_0, x_1)$ is extended with a dummy attribute $x_0 = 1$.

**Question 2:** $\boxed{\text{1 point}}$ The three most common problems in supervised learning are classification, regression, and logistic regression. For each problem, give a practical example, including a description of the input and output.

Name: ...........................................................................

**Question 3:** ⟨1 point⟩ What is the role of the error measure in supervised learning? How can the perceptron learning algorithm be adapted to different error measures for classification?

**Question 4:** ⟨1 point⟩ $V$-fold cross-validation is a popular technique for model selection. Approximately how much higher is the running time compared to regular validation?

**Question 5:** ⟨1 point⟩ Gradient descent is frequently used to find the minimum of a convex function. How does the effectiveness of gradient descent depend on the learning rate $\eta$?

**Question 6:** ⟨1 point⟩ What is the kernel trick? Explain why the kernel trick is often combined with support vector machines to learn a classifier.

Name: ................................................................

**Question 7:** [1 point] The Policy Iteration algorithm for Markov decision processes iteratively computes a sequence of policies $\pi_1, \pi_2, \ldots, \pi_k$ and a sequence of value functions $V_1, V_2, \ldots, V_k$. In a given iteration $k$ of the algorithm, what is the relation between $\pi_k$ and $V_{k+1}$?

**Question 8:** [1 point] Describe the TD(0) algorithm for policy evaluation in a discounted Markov decision process. Name one advantage and one disadvantage of TD(0) over the Least-Squares temporal difference (LSTD) learning algorithm.

**Question 9:** [2 point] Consider the Explore-then-commit (ETC) algorithm for a two-armed bandit problem with mean rewards $\mu_1$ and $\mu_2 = \mu_1 - \Delta$, with rewards bounded in $[0, 1]$. Let $\widehat{\mu}_1$ and $\widehat{\mu}_2$ be the empirical mean rewards obtained from the two arms in the exploration phase of ETC, respectively. Give a lower bound on the sample size $m$ necessary for guaranteeing $\widehat{\mu}_1 \geq \widehat{\mu}_2$ with probability at least $1 - \delta$!
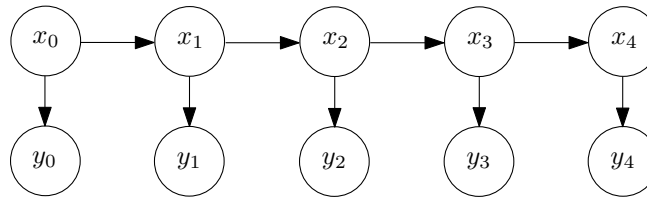**Hint:** Use Hoeffding's inequality that guarantees

$$\frac{1}{m}\sum_{t=1}^{m} X_t - \mathbb{E}\left[X_1\right] \leq \sqrt{\frac{\log(1/\delta)}{m}}$$

with probability at least $1 - \delta$ if $X_t \in [0, 1]$ for all $t$.

**Question 10:** 2 point

- Consider the Bayesian network depicted below. Write down the corresponding factor graph using squared nodes to denote factors and circle nodes to denote variables.



- Assuming that $x_1$ can take three discrete values and that $x_3$ is a binary variable, explain a possible way to compute the marginal $p(x_1, x_3)$ using the least possible number of runs of the belief propagation algorithm. You don't need to write down the messages, just state, for each run, which variables do you need to clamp, and what is the result you want to obtain.

Name: ...........................................................................

**Question 11:** ⟨2 point⟩

- Continuous random variables $x$ come independently from a probability distribution. According to a model $\mathcal{H}_1$, this probability distribution is

$$P(x|m, \mathcal{H}_1) = \frac{1}{2}(1 + mx), \qquad x \in (-1, +1), \qquad (1)$$

where $m$ is the only model parameter, with value between $-1$ and $+1$.

After observing three data points of $x$, $D = \{\frac{1}{3}, \frac{1}{2}, \frac{3}{5}\}$, find $m$ analytically (assume a uniform prior for $m$).

- A simpler explanation $\mathcal{H}_0$ for observing $D$ is that those variables come from a uniform probability distribution

$$P(x|m, \mathcal{H}_0) = \frac{1}{2}. \qquad (2)$$

Given the data $D$, what is the evidence for $\mathcal{H}_0$? What is the evidence for $\mathcal{H}_1$?

*Remember that* $\int_a^b m^n dm = \left( \frac{m^{n+1}}{n+1} + C \right)\Big|_a^b$

The total score is: _____ /15.