SOLUTIONS

CS689: Machine Learning - Fall 2019

Final Exam

Dec 13, 2019	Name:	
--------------	-------	--

Instructions: Write your name on the exam sheet. The duration of this exam is two hours. No electronic devices may be used during the exam. You may consult either your notes or the Murphy text during the exam (you must choose one). Up to a three inch stack of notes is permitted. No other sources are permitted. Sharing of notes/texts during the exam is strictly prohibited. Show your work for all derivation questions. Provide answers that are as detailed as possible for explanation questions. Attempt all problems. Partial credit may be given for incorrect or incomplete answers. If you need extra space for answers, write on the back of the preceding page. If you have questions at any time, raise your hand.

Problem	Topic Page		Points	Score
1	Loss Functions	1	10	
2	Maximum Likelihood	2	10	
3	Generalized Linear Regression	3	10	
4	Logistic Regression	4	10	
5	Lagrangian Duality	5	10	
6	Mixture Models	6	10	
7	Factor Analysis and Autoencoders	7	10	
8	Bayesian Inference	8	10	
9	Experiment Design	9	10	
10	Applications	lications 10 10		
Total:			100	

- 1. (10 points) Loss Functions Suppose we have a classification problem where $x \in \mathcal{X}$ and $y \in \mathcal{Y} = \{1, ..., C\}$. Suppose that f is a prediction function satisfying $f : \mathcal{X} \to \mathcal{Y}$. Under what circumstances would it be sensible to assess the performance of f using the squared loss? When would it not be sensible? Explain your answers.
- * For a classification problem, we generally care about the classification error or zero-one loss.
- * In the case where C=Z, the possible class
 labels are {1, 2} and the resulting values of
 squared loss and zero-one loss are given below.

У	f(x)) + f(x)	(1- f(x1)
	1	0	0
1	Z	1	1
2	2	0	0

- or C=2 and using squared loss is sensible.
- * For C>2 the squared loss is larger when the difference in label values is larger, even though the label values are not meaningful. This is not sensible.

- 2. (10 points) Maximum Likelihood: Suppose we have a binary probabilistic classifier $P_{\theta}(Y = y|X = x)$. Given a data set of feature vectors $\mathcal{D} = \{x_n | 1 \leq n \leq N\}$, we are interested in modeling the distribution of the classifier output probabilities $\pi_n = P_{\theta}(Y = 1|X = x_n)$ using a Beta distribution $P(\Pi = \pi|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\pi^{a-1}(1-\pi)^{b-1}$.
- (a) (5 points) What is the log likelihood of the model P(Π = π|a, b) given the observations π₁, ..., π_N?

(b) (5 points) The digamma function $\psi(x)$ is a special function $\psi(x) = \frac{d \log \Gamma(x)}{dx}$ that provides the derivative of the log of the gamma function $\Gamma(x)$. Use this fact to derive the gradient vector for the model $P(\Pi = \pi | a, b)$. Show your work.

- 3. (10 points) Generalized Linear Regression: Consider the probabilistic linear regression model $P(Y = y | \mathbf{X} = \mathbf{x}) = \mathcal{N}(y; \mathbf{w}^T \mathbf{x} + b, \sigma^2)$. This model asserts that the variance is constant for all \mathbf{x} . Suppose we are interested in generalizing the model so that the variance is a function of \mathbf{x} as well.
- (a) (5 points) Describe how you could modify the model to also make the variance datadependent. Provide an updated form for P(Y = y|X = x).

(b) (5 points) Suppose we only care about the mean squared error that the learned model you described in part (a) achieves when making predictions on test data. Do we need the learned data-dependent variance component of the model when making predictions? Explain your answer.

No. To make a prediction, we typically output the mean of the predictive distribution Ep(VIX)[Y]. For a normal model this is just a. For the linear gaussian model it is wtx+b, which does not depend on the variance at all.

4. (10 points) Logistic Regression Equivalence: The binary logistic regression model with labels $y \in \{-1, 1\}$ and features $\mathbf{x} \in \mathbb{R}^D$ can be written as

$$P(Y = y|\mathbf{X} = \mathbf{x}) = \frac{1}{1 + \exp(-y(\mathbf{w}^T\mathbf{x} + b))}$$

using one weight vector $\mathbf{w} \in \mathbb{R}^D$ and one bias b. On the other hand, the multi-class logistic regression model with $y \in \{1, 2, ..., C\}$ can be written as

$$P(Y = y | \mathbf{X} = \mathbf{x}) = \frac{\exp(\mathbf{w}_y^T \mathbf{x} + b_y)}{\sum_{y'=1}^{C} \exp(\mathbf{w}_{y'}^T \mathbf{x} + b_{y'})}$$

using one weight vector \mathbf{w}_y and one bias b_y per class. Despite their apparent differences, show that the two models are exactly equivalent in the case where C=2.

We will show that the multi-category case with C=2 reduces to the binary case.

We divide the top and bottom by exp(wz*x+bz) to get.

Non let w= wz-w, b= bz-b, we hare:

If we map y=1 to \$\hat{y}=1 and y=2 to \$\hat{y}=1, we have $P(Y=y|X=x)=P(\hat{Y}=\hat{y}|X=x)=\frac{1}{1+exp(-\hat{y}(\hat{w}^Tx+\hat{b}))}$

5. (10 points) Lagrangian Duality Consider the alternate optimization problem for linear regression given below:

$$\arg \min_{\mathbf{w}, \eta} \sum_{n=1}^{N} \eta_n^2$$
s.t. $\forall n, \quad \eta_n = y_n - (\mathbf{w}^T \mathbf{x}_n)$

$$\|\mathbf{w}\|_2^2 \leq C$$

(a) (5 points) Write down the Lagrangian function for this optimization problem.

(b) (5 points) Derive the Lagrange dual for this optimization problem. Show your work.

$$\frac{\partial \mathcal{L}(W,3,\lambda,\delta)}{\partial \beta_n} = \lambda \beta_n - \lambda_n = \beta \beta_n = \frac{\lambda_n}{2}$$

$$q(\lambda, \delta) = \sum_{n=1}^{\infty} \frac{\lambda_{n}^{2}}{4} - \sum_{n=1}^{\infty} \lambda_{n} \left(\frac{\lambda_{n}^{2}}{2} - (y_{n} - (\frac{1}{2\delta} \sum_{n=1}^{\infty} \lambda_{n}^{2} \times \lambda_{n}^{2})^{T} \times_{n}) \right)$$

$$- \chi \left(C - \sum_{n=1}^{\infty} \sum_{n=1}^{\infty} \lambda_{n} \lambda_{n}^{2} \times_{n}^{T} \times_{n} \right)$$

- 6. (10 points) Mixture Models Suppose we have a data set $\mathcal{D} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$ consisting of data vectors $\mathbf{x}_n = [x_{1n}, ..., x_{Dn}]$ where each data dimension represents an integer count.
- (a) (5 points) Describe how we could model the distribution P(X = x) as a mixture model. Be specific and provide equations to support your description.

We can use a mixture of products of count distributions such as geometric or poissin.

We have &

$$b(x=x) = \sum_{K=1}^{K=1} b(x=x|S=K)b(S=K) = \sum_{K=1}^{K=1} \Theta^{K} \cdot \prod_{D} (1-\prod^{SK})_{x} \prod^{SH}$$

(b) (5 points) Suppose we want to learn the model using direct numerical optimization. Provide the specific objective function you would use and explain your answer.

Since the 2 variables are not observed, we must maximize the log marginal likelihood:

7. (10 points) Factor Analysis and Autoencoders

(a) (5 points) Consider the factor analysis model is given by P(Z = z) = N(z; 0, I) and P(X = x|Z = z) = N(x; w^Tz, Ψ) where Ψ is a positive definite matrix. What additional restrictions are needed on Ψ for this model to learn a useful latent representation for the data? Explain your answer.

We need to restrict 4 to be a diagonal matrix.

The reason is that this will force the structure in the data to be explained by the latent factors 2.

Otherwise, 4 itself can model the covariance in the data. Setting 4 = 5°I yields P(A and will also result in learning useful structure.

(b) (5 points) Explain why there is no modeling benefit to using a linear autoencoder with more than one hidden layer.

* If we add an extra set of layers we have:

hi= WiTX; hz= WzTh,; hz= VzThz; r= V,Thz



* Since the layers are linear, we have s hz = Wz W, x and r= V, Vz hz = vhz

* Thus, adding exten layer does not increase representational capacity.

8. (10 points) Bayesian Inference: Suppose that $x \in \mathbb{N}$ is a count random variable that follows the negative Binomial distribution:

$$P(X = x | \theta, r) = {x + r - 1 \choose x} \cdot (1 - \theta)^r \theta^x$$

We use a Beta prior distribution as defined below. Use this information to answer the following questions.

$$P(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha - 1} (1 - \theta)^{\beta - 1}$$

(a) (5 points) Given data $\mathcal{D} = \{x_n | n = 1, ..., N\}$, what is the likelihood of \mathcal{D} given θ ? Explain your answer and simplify as much as possible.

$$+ P(D|\Theta) = \prod_{n=1}^{N} P(X:x \cap |\Theta, \Gamma) = \prod_{n=1}^{N} {\binom{x_n \cdot \Gamma - 1}{x_n}} (1 - \Theta)^{\Gamma} \Theta^{x_n}$$

$$= \prod_{n=1}^{N} {\binom{x_n \cdot \Gamma - 1}{x_n}} \cdot (1 - \Theta)^{N\Gamma} \cdot \Theta^{\sum_{n=1}^{N} x_n}$$

(b) (5 points) Given data D = {x_n|n = 1,...,N}, what is the posterior distribution of θ given D? Show your work and explain your answer.

* The posterior distribution is proportional to the likelihood times the prior. We expand, drop constants and identify the form of posterior.

- 9. (10 points) Experiment Design: Suppose we have a data set containing labeled instances $(\mathbf{x}_{tn}, y_{tn})$ for a collection of individuals n = 1, ..., N and time points t = 1, ..., T.
- (a) (5 points) If we want to assess how well a supervised model can perform when applied to a future set of time points for the same individuals, what experimental design should we use? Explain your answer.

We first need to train the model. We then need to evaluate it. The most basic design is train-test.

Since we want to evaluate generalization to future times, we need to simulate this using the data we have. We can pick a value To < T and as the training set use Dtr = {(X+n, Y+n)| +>To, 1≤n ∈ N}. As the test set we use Dte = {(X+n, Y+n)| +>To, 1≤n ∈ N}

(b) (5 points) Suppose the data set contains many features and we start by using all available data vectors \mathbf{x}_{tn} to learn and apply a dimensionality reduction model, obtaining lower-dimensional feature vectors \mathbf{x}'_{tn} . We then apply the procedure you described in part (a) to a data set of $(\mathbf{x}'_{tn}, y_{tn})$ instances. Is the end-to-end experiment design still correctly assessing how well this approach can perform when applied to data from new time points? Explain your answer.

No. The supervised model depends on the values x'tm, and x'tm are inferred using the complete data set. It is thus using information from the test set during training. If the distribution of x changes through time, this procedure would underestimate the error.

- 10. (10 points) Applications: Suppose we have a data set of completely observed real-valued vectors $\mathcal{D} = \{\mathbf{x}_n | n=1,...,N\}$ where $\mathbf{x}_n \in \mathbb{R}^D$. Our goal is to learn a dimensionality reduction model for these data. However, at deployment time, we expect to encounter incomplete data vectors (vectors with missing entries) and we must still be able to produce reduced dimensional representations for them. Describe a solution for modeling these data that does not require an auxiliary imputation model. Your answer should cover (1) a description of your model, (2) your approach to learning the model, and (3) how the model is applied to incomplete data cases. Provide supporting equations and justify your choices.
- We can use a factor analysis midel for this

 task. The model is P(z): N(z; o, I), P(x|z): N(x; w⁷z, 4)

 where 4 is a diagonal matrix and w is the

 factor loading matrix. We make this choice because

 the model can directly support missing data.

 The model can directly support missing data.

 We can learn the model as usual from

 complete data using direct marginal likelihood
- 3 Given an incomplete data instance x, we need to compute P(Z=Z|X°=x°). We can get this using:

maximization (marginalizing over Z).