

# Midterm One, CS-UY 4563

---

## Basic Info

**Date:** Monday 3/9 from 9am - 10:15am **Location:** Our usual classroom (JABS 775B) **Supplies needed:** Pencil or pen. I will provide scrap paper. **Supplies allowed:** One-double sided sheet of paper containing whatever you like. Equations, notes, examples, etc. This can be handwritten or typewritten -- there are no restrictions.

## Preparation

The best way to prepare is to review your homework and labs. If there are any homework problems you did not understand or fully solve, please ask about them in the review session on Wednesday, or ask an instructor during office hours!

## Question Types

- Short answers similar to the homework questions.
- True/false, with a short statement to justify your answer.
- "Always, sometimes, never" questions. I give you a statement, you say if it's always true, sometimes true, or never true, and give a short written justification.
- Data diagnostics: I show you a plot or table of data and ask you what machine learning method, or data modification, might be appropriate for the data.

## Topics

---

Below are a list of things you should know, roughly catagorized into topics.

### Introduction to Machine Learning

- What is *supervised* learning?
- What is *classification*? What is *regression*?
- What is a *predictor variable* what is a *target variable*? What is a *data matrix*?
- Loss minimization framework: what is a *model*, the *model parameters*, and the *loss function*. The objective of supervised learning in the loss minimization framework is to choose model parameters which minimize the loss.

### Simple Linear Regression

- What is the model for simple linear regression? What are the parameters?
- Know these different loss functions for regression:  $\ell_2$ ,  $\ell_1$ ,  $\ell_\infty$ . You should be able to come up with other reasonable loss functions on your own.
- Why might you use one loss over another? Understand the intuition for why  $\ell_2$  and  $\ell_1$  are more

robust to outliers than  $\ell_\infty$ . Come up with an example problem where you might want to use  $\ell_\infty$  loss any way.

- To minimize a function with multiple inputs, *set all partial derivatives to 0* to find extreme points.
- Or run brute force search!
- How do you set  $\beta_0$  and  $\beta_1$  to minimize the squared loss for simple linear regression. *You should be able to derive these expressions entirely on your own.*
- What are some reasons we might get a poor fit during linear regression? Which problems can we address? Which are inherent issues with the data or model?
- Given a non-linear problem, know how to transform variables to convert it into a linear model (when possible).

## Multiple Linear Regression

- What is the multiple linear regression model? Write it using matrix/vector notation.
- What losses could be used for multiple linear regression? Write using matrix/vector notation. E.g. the squared loss would be  $L(\vec{\beta}) = \|X\vec{\beta} - \vec{y}\|_2^2$ .
- Adding a column of all 1s on to the front of the data matrix allow us to eliminate the intercept parameter. Adding a column of all 2s would have worked as well. Or a column of all -2s. We would get the exact sample predictions from a linear regression model. Why? A column of all 0s would not work. Why?
- What is the *gradient* of a function? Setting the gradient to  $\vec{0}$  is the same as setting all partial derivatives to 0.
- What is the gradient for multiple linear regression under the squared loss? Be able to compute the gradient for other loss function (as in HW 2, Prob. 3) using the multivariate chain rule.
- Make sure you sanity check all gradient computations: if  $L$  is a function which takes a  $d$ -dimensional vector as input and outputs a number (the loss),  $\nabla L$  takes in a  $d$ -dimensional vector and outputs a  $d$ -dimensional vector (the gradient).
- Be able to compute gradients of other simple multivariate functions. E.g. what is the gradient of  $f(x) = \langle a, x \rangle$ , or  $f(x) = e^{\langle a, x \rangle}$ ?
- Once you have the gradient, how do you solve for the optimal parameter vectors  $\vec{\beta}$  for linear regression under squared loss: Compute  $(X^T X)^{-1} X^T y$ .
- Data scaling and centering has *no effect* on the optimal parameters for multiple linear regression. Many students missed this question on HW2.
- How can multiple linear regression be used with data transformations to fit a  $q$ -degree polynomial? A sinusoidal curve? Other models?
- What is *one-hot encoding*? Why is it necessary? See beginning of lecture 4.

## Model Selection, Train-Test Split

- What is *overfitting*? What is *generalization*? You don't need to know mathematical definitions, but understand what they mean intuitively.
- Does better loss on the training set mean better generalization? Always, sometimes, or never?
- What is *model selection* formally? See Slide 5 in Lecture 5.
- Here are some examples of model classes of varying complexity: polynomials with different degrees, bag-of-words with different values of  $n$  in the  $n$ -gram computation, time series models

with different delays (like on the lab). You might use model selection in all of these cases to determine the right model complexity.

- How to split a data into a training set and a test set.
- What is the definition of *population risk*? Our objective in machine learning is almost always to find models that minimize population risk.
- The expected value of the *empirical risk on the test set* is equal to the population risk. Why is this observation useful?
- What is  $k$ -fold cross validation? Why would you use it instead of train-test split? Why would you not use it?

## Feature Selection and Regularization

- What is regularization? How does it change the loss-minimization framework?
- For what types of machine learning problems do you need to use feature selection or regularization for?
- Typically we need one or the other for over-parameterized models, where the number of predictor features is larger than the number of training data points. Why are you essentially *guaranteed* to overfit in this setting when you have a linear regression model?
- Both approaches can also prevent overfitting even when you aren't obviously overparameterized.
- What is *ridge regularization* ( $\ell_2$ ) and *LASSO regularization* ( $\ell_1$ )?
- If you add a regularization penalty, do the parameters in your optimal  $\vec{\beta}$  increase or decrease in comparison to if you had minimized the loss without regularization?
- Note that not *every* parameter necessarily decreases, but we proved that if you add a regularization term of  $\lambda \|\vec{\beta}\|$  with  $\lambda > 0$ , then  $\|\vec{\beta}\|$  decreases. This is true for *any regularization norm*.
- What's one advantage of feature selection over regularization? It leads to a simpler model, where  $y$  is predicted using fewer parameters.
- How is regularization justified from a Bayesian perspective?
- What is the maximum a posterior (MAP) estimate for  $\vec{\beta}$  under an assumption that our model is linear with Gaussian noise? What about under Laplace noise?
- Be able to derive why a Gaussian prior on the parameters  $\vec{\beta}$  tells us to use  $\ell_2$  regularization and why a Laplace prior tells us to use  $\ell_1$  (LASSO regularization). What is the *log likelihood function*?

## Naive Bayes + Probabilistic Modeling

- Be able to come up with a probabilistic model for a given data set and understand how you would learn parameters of the model from data. Homework 3, Problem 1 is good practice! Talk to us if you're struggling.
- Definition of *joint probability*  $p(a, b)$  and *conditional probability*  $p(a | b)$ . Which is larger?
- Bayes Rule:  $p(a | b) = \frac{p(b|a)p(a)}{p(b)}$ .  $p(b | a)$  is called the *likelihood*,  $p(a)$  is the *prior*,  $p(b)$  is the *evidence* and  $p(a | b)$  is the *posterior*.
- What is a *bag-of-words vector representation* of a document? Why is it useful?
- What probabilistic model for binary vectors justifies the *Naive Bayes Classifier*? How are the parameters of this model learning from past data.

- What is the final *Naive Bayes* classification rule?
- What is the equation for a Gaussian/normal PDF? For a Laplace PDF?
- What is a *maximum a posteriori (MAP)* estimate for predicting a label or regression value  $y$  given data  $\vec{x}$  with parameters  $\vec{\theta}$ ? Choose  $y$  to maximize  $p(y \mid \vec{x}, \vec{\theta})$ .
- What is a *maximum a posteriori* estimate for determining model parameters? An example of this is when we used Bayes rule to choose an optimal  $\beta$  for linear regression in Lecture 8, under a specific probabilistic model.

## Linear Classifiers and Logistic Regression

- What is a linear classifier? How would you right down a model for linear classification?
- What is the 0-1 loss?
- The logistic function  $h_{\vec{\beta}}(\vec{x}) = \frac{1}{1+e^{-\langle \vec{\beta}, \vec{x} \rangle}}$ .
- What is the model and loss function which define *logistic regression*.
- What is the *error rate*, *precision*, and *recall* of a classifier? How would you adjust a classifier to improve precision or recall? In what real world scenario would you care about precision more? Recall more?
- How do the *one-vs-all* and *one-vs-one* methods work for multiclass classification? In general, one-vs.-one leads to more accurate results, but at the cost of increase computational complexity. What is that complexity overhead?
- What is a *confusion matrix*.

## Optimization and Gradient Descent

- Why do we need to use search methods to minimize loss functions instead of just figuring out when the gradient equals  $\vec{0}$ .
- How would you do a brute force search for the optimal parameters  $\vec{\beta}$  of a model? Why would this be slow when  $\vec{\beta}$  has more than a few dimensions?
- The gradient descent update is  $\vec{\beta} \leftarrow \vec{\beta} - \eta \nabla L(\vec{\beta})$  where  $\eta$  is the *stepsize*, aka *learning rate*, parameter.
- $\lim_{\eta \rightarrow 0} L(\vec{\beta} + \vec{v}) - L(\vec{\beta}) = \eta \langle \nabla L(\vec{\beta}), \vec{v} \rangle$ .
- Why is the gradient descent update  $\vec{\beta} \leftarrow \vec{\beta} - \eta \nabla L(\vec{\beta})$  guaranteed to decrease  $L(\vec{\beta})$  for a sufficiently small step-size  $\eta$ ?
- Why is gradient descent sometimes called *steepest descent*?
- **You don't need to know anything past slide 25 in Lecture 10.**