

Homework 4: Underdetermined Problems and Model Validation

Due: 10/31/19

1. *Linear dependence.* Prove or disprove (by providing a counterexample) each of the following statements.

- (a) If the rows of a matrix X are linearly dependent, then X is not invertible.
- (b) If a matrix X has full column rank, then $X^T X$ is invertible.
- (c) If X is an $m \times n$ matrix with full row rank but *fat* ($n > m$), then $X^T X$ is not invertible.
- (d) Why is part (c) problematic when we are solving a least-squares problem? What does it mean (in words) if our data matrix is fat?

2. *Rats!*

Your friend Alice works in a lab where she studies the genetic indicators of rat intelligence. To do this, she sequences their DNA, and tests how long it takes for the rats to finish a maze. She has collected all of the rats' genetic data into a matrix X . Each row contains data for one rat and the columns are different alleles (a variant of a gene). The entry $X[i, j]$ is 1 if rat i has allele j but 0 otherwise. She's collected all of the rat maze finish times into the vector y .

One day she tells you that she has run a linear regression on her dataset to understand which alleles are important for rat intelligence. She has found a model w_A that fits the problem with very small residual sum of squares error.

- (a) Alice is excited because she's found two alleles with very large coefficients which she thinks might be important for understanding rat intelligence. The coefficient for the first allele, $w_A[1]$, is very large and positive. The coefficient for the second allele, $w_A[2]$, is very negative. How do you interpret these coefficients? What would you predict about how long it would take a rat with each of these genes to finish the maze?
- (b) Bob works in the same lab as Alice. The next day, he tells you that he's solved the same least squares problem with an error equally as small. But his coefficients,

w_B , are very different from Alice's! His coefficient for the first allele, $w_B[1]$, is very negative while $w_B[2]$ is very positive. Why do you think Bob's coefficients are so different from Alice's?

- (c) Candace, who also works in the same lab, tells you that she has found a rat with a genetic code that has never been seen before. All other lab rats had both allele 1 and allele 2. This new rat only has allele 1. Alice's model predicts that this rat will take 535 years to finish the maze. Bob's model predicts the rat will finish in -534 years. How will you answer Alice, Bob and Candace if they ask you to make a prediction?
- (d) If Bob and Alice gave you their dataset X and y , what would you do to make a better prediction? *(There isn't necessarily a right answer here.)*

3. Ridge regression and the SVD

Recall the singular value decomposition from class. The SVD decomposes a matrix $X \in \mathbf{R}^{n \times d}$ into a product of simpler matrices

$$X = U\Sigma V^T,$$

where

- $U \in \mathbf{R}^{n \times r}$ and $V \in \mathbf{R}^{d \times r}$ have orthonormal columns

$$U^T U = I, \quad V^T V = I.$$

- $\Sigma \in \mathbf{R}^{r \times r}$ is diagonal with $\Sigma_{ii} = \sigma_i > 0$ for $i = 1, \dots, r$.
- $r \leq \min(d, n)$ is the rank of X .

(This version is sometimes called the thin SVD.) In class we showed how to use the SVD to solve the least squares problem. In particular, we showed that $w^{\text{SVD}} = V\Sigma^{-1}U^T y$ satisfies the normal equations.

In class, we saw that when X is not full rank, least squares has no unique solution. We introduced ridge regression to deal with these underdetermined problems. In this problem we will draw a connection between w^{SVD} and w^{ridge} .

We will show that the w^{SVD} is the solution to least squares that minimizes the Euclidean norm $\|w\|$. That is, if we have a vector w s.t. $\|y - Xw\| = \|y - Xw^{\text{SVD}}\|$, then $\|w\| \geq \|w^{\text{SVD}}\|$.

- (a) Let w be some solution to the least squares problem. It must satisfy the normal equations. Rewrite the normal equations ($X^T X w = X^T y$) using the matrices U , V , and Σ .

- (b) Define a matrix V^c whose columns form an orthonormal basis for the orthogonal complement of V . That is, $(V^c)^T V = 0$ and $(V^c)^T V^c = I_{d-r}$.

We can decompose $w = Vw^\parallel + V^c w^\perp$ where $w^\parallel = V^T w$ and $w^\perp = (V^c)^T w$. Use the normal equations to derive a formula for w^\parallel in terms of U , Σ , and V . In particular, this shows that any solution w to least squares must have the same w^\parallel .

- (c) Decompose $w^{\text{SVD}} = Vw^{\text{SVD}\parallel} + V^c w^{\text{SVD}\perp}$ using the same decomposition. Derive a formula for $w^{\text{SVD}\parallel}$ and $w^{\text{SVD}\perp}$.
- (d) Prove that $\|w\|^2 = \|Vw\|^2$ if $V^T V = I$.
- (e) Write out the Euclidean norm of w in terms of w^\parallel and w^\perp . Conclude that for any solution w to least squares, $\|w\| \geq \|w^{\text{SVD}}\|$
Hint: Use the Pythagorean Theorem.

Through parts (a) - (d), we have shown that w^{SVD} is the solution to least squares with minimum norm. Recall that the ridge regression estimator,

$$w^{\text{ridge}}(\lambda) = \operatorname{argmin} \|y - Xw\|^2 + \lambda\|w\|^2,$$

minimizes the mean squared error and a weighted norm of w .

- (f) Show that

$$\lim_{\lambda \rightarrow 0} w^{\text{ridge}}(\lambda) = w^{\text{SVD}}$$

Hint: Write w^{ridge} and w^{SVD} in terms of the σ_i 's.

4. Plotting bias and variance

In this problem, we'll investigate the bias and the variance of two different estimators: a linear estimator and a cubic estimator. We'll see (once again) that fitting the data more precisely is not always a good idea.

Work through the notebook <https://github.com/ORIE4741/homework/blob/master/BiasVariance.ipynb>.

- (a) Suppose we have a sinusoid function $f(x) = 10 \sin(x)$. Our dataset \mathcal{D} will consist of seven datapoints drawn from the following probabilistic model.
 For each datapoint we randomly draw x_i uniformly in $[0, 6]$ and observe a noisy

$$y_i = f(x_i) + \epsilon_i,$$

where ϵ_i is some noise drawn from a standard normal distribution $\mathcal{N}(0, 1)$.

Generate a sample dataset from this distribution. Plot this dataset \mathcal{D} and the true function $f(x)$.

- (b) Fit a linear model $l(x) = w_0 + w_1 x$ to \mathcal{D} .
 Plot this linear model $l(x)$ together with \mathcal{D} and $f(x)$.

- (c) Fit a cubic model $c(x) = w_0 + w_1x + w_2x^2 + w_3x^3$ to \mathcal{D} . Plot this cubic model $c(x)$ together with D and $f(x)$.
- (d) Repeat parts (b) and (c) for 1,000 different randomly drawn sets \mathcal{D} . Average the 1,000 linear models you generated to get the average linear model $\bar{l}(x)$. Plot $\bar{l}(x)$ with $f(x)$.
Generate the average cubic model, $\bar{c}(x)$, in the same way and plot it together with $f(x)$.

We can decompose the expected out-of-sample error into the squared bias,

$$\mathbb{E}[\text{bias}^2(x)] = \mathbb{E}_x[(f(x) - \bar{g}(x))^2],$$

and variance,

$$\text{Var}(x) = \mathbb{E}_D[(g^D(x) - \bar{g}(x))^2],$$

where $g^D(x)$ is the model generated with dataset \mathcal{D} .

- (e) Compute the squared bias of $\bar{l}(x)$ and $\bar{c}(x)$. Which model has smaller squared bias?
- (f) Compute the variance of $\bar{l}(x)$ and $\bar{c}(x)$. Which model has smaller variance? How do you interpret this? Which model has smaller overall error?
- (g) How do you think your results would depend on the number of points in the data set \mathcal{D} ? Feel free to perform an experiment to check. How many points would you need before the *opposite* model (from your answer in the previous question) has smaller overall error?
- (h) Instead of sampling new data to compute the bias and variance of our model, we could use a bootstrap estimator to get more use out of the few data points we have. Try this for a few different data set sizes and report on your results. How big a data set is needed for the bootstrap to give a reliable estimate of the bias and variance?