Chaper 5: Matrix Approach to Simple Linear Regression

```
Matrix:
A m by n matrix B is a grid of numbers with m rows and n columns.
B = b_{11}... b_{1n}
             •
       •
       •
      b<sub>m1</sub>... b<sub>mn</sub>
Element b_{ik} is from the ith row and kth column.
A vector b is a matrix with 1 column:
b = b_1
       •
       •
       •
      \mathbf{b}_{n}
A transpose of a m by n matrix B is a n by m matrix B'
B′ =
             b<sub>11</sub>... b<sub>m1</sub>
       •
             •
             •
       •
      b<sub>1n</sub>... b<sub>mn</sub>
A product of a m by n matrix (B) and a n by p matrix (A) is:
BA =
          (b_{11}a_{11}+...+b_{1n}a_{n1}) ... (b_{11}a_{1p}+...+b_{1n}a_{np})
                           •
       •
       (b_{m1}a_{11}+...+b_{mn}a_{n1}) ... (b_{m1}a_{1p}+...+b_{mn}a_{np})
which has dimension m by p. (See board for detail.)
```

```
An n by n (square) matrix I is called the identity matrix if
I = 1 0 \dots 0
      0 1 ...
                       0
                 1 0
      0 ... 0 1 (We say "1 on the diagonal")
Note that Ib = b (assuming the dimensions line up).
Some square matrices have an inverse: B^{-1} which is matrix such that BB^{-1} = I.
A regression model can be expressed as:
\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} where
\mathbf{y} = \mathbf{y}_1 \mathbf{x} = 1 \mathbf{x}_1 \boldsymbol{\beta} = \boldsymbol{\beta}_0 \boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}_1
                        • β<sub>1</sub>
       •
                   •
                   •
                        •
               • •
      •
            1 x,
      y<sub>n</sub>
                                                   \mathcal{E}_n
The matrix x is called the design matrix and \beta are the coefficients.
e is a random vector. It has mean 0 =
                                                   0
                                                    0
and covariance matrix: cov(e)=
\sigma^2
      0
                          0
            ...
                   ...
      σ<sup>2</sup> ... ...
0
                          0
                 \sigma^2 0
```

0 ... 0  $\sigma^2$  The ikth entry is  $cov(\varepsilon_i, \varepsilon_k)$ .

A nice thing about matrices is that the regression coefficient estimates now have a simple form:

 $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  where  $\mathbf{b} = (\mathbf{b}_0 \ \mathbf{b}_1)'$ . (Derive on board.)

**Y-hat** =  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  = **Hy** and **H** is called the "hat" matrix.

 $cov(\mathbf{b}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$  which is estimated by  $MSE(\mathbf{X}'\mathbf{X})^{-1}$ .

```
Chapter 6: Multiple Regression I
```

```
In the simple (i.e. only one covariate) linear regression

y = beta0 + beta1 x + e, the model was motivated by fitting a line through a scatterplot.

A multiple linear regression models the mean of y as a linear function of more than one

covariate:

y = beta0 + beta1 x1 + ... + betap xp + e

with e independent, normal(0,sigma^2)

Note that the regression model is:

E(y) = beta0 + beta1 x1 + ... + betap xp

or (write matrix formulation)

These are first order models since they do not include products of xs.

(i.e. E(y) = beta0 + beta1 x1 + beta2 x2 + beta3 (x1*x2) is a second order model.)

Polynomial models are multiple regressions:

<math>E(y) = beta0 + beta1 x + beta2 x^2
```

When p = 2, the betas are estimated by fitting a plane through a cloud of data (see figure). For p>2, it's a higher dimensional generalization.

(write matrix formulation for estimators)

## Example: Cheddar Cheese

**Reference:** Moore, David S., and George P. McCabe (1989). Introduction to the Practice of Statistics. Authorization:

Description: As cheese ages, various chemical processes take place that determine the taste of the final product. This dataset contains concentrations of various chemicals in 30 samples of mature cheddar cheese, and a subjective measure of taste for each sample. The variables "Acetic" and "H2S" are the natural logarithm of the concentration of acetic acid and hydrogen sulfide respectively. The variable "Lactic" has not been transformed. Variable Names:

- 1. Case: Sample number
- 2. Taste: Subjective taste test score, obtained by combining the scores of several tasters
- 3. Acetic: Natural log of concentration of acetic acid
- 4. H2S: Natural log of concentration of hydrogen sulfide
- 5. Lactic: Concentration of lactic acid

Our goal is to predict Taste.

First model: Taste; = beta0 + beta1 Acetic; + beta2 H2S; + beta3 Lactic; + e; Call: lm(formula = taste ~ Acetic + H2S + Lactic, data = cheese) Coefficients: Estimate Std. Error t value Pr(>|t|)(Intercept) -28.8768 19.7354 -1.463 0.15540 Acetic 0.3277 4.4598 0.073 0.94198 H2S 3.9118 1.2484 3.133 0.00425 \*\* 19.6705 8.6291 2.280 0.03108 \* Lactic Residual standard error: 10.13 on 26 degrees of freedom Multiple R-Squared: 0.6518, Adjusted R-squared: 0.6116

F-statistic: 16.22 on 3 and 26 DF, p-value: 3.81e-06

Go through what the various things in the table above mean.

Do parameter estimates change depending on what else is in the model?

```
Yes! (when the covariates are correlated themselves)
Call:
lm(formula = taste ~ Acetic, data = cheese)
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) -61.499
                        24.846 -2.475 0.01964 *
Acetic
             15.648
                         4.496 3.481 0.00166 **
Residual standard error: 13.82 on 28 degrees of freedom
Multiple R-Squared: 0.302, Adjusted R-squared: 0.2771
F-statistic: 12.11 on 1 and 28 DF, p-value: 0.001658
Taste, = beta0 + beta1 Acetic, + beta2 H2S, + beta3 Lactic, + e,
Call:
lm(formula = taste ~ Acetic + H2S + Lactic, data = cheese)
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) -28.8768
                       19.7354 -1.463 0.15540
Acetic
             0.3277
                        4.4598 0.073 0.94198
H2S
             3.9118
                       1.2484 3.133 0.00425 **
                        8.6291
                                 2.280 0.03108 *
Lactic
            19.6705
Residual standard error: 10.13 on 26 degrees of freedom
Multiple R-Squared: 0.6518, Adjusted R-squared: 0.6116
F-statistic: 16.22 on 3 and 26 DF, p-value: 3.81e-06
```



Analysis of Variance Table

Response: taste Df Sum Sq Mean Sq Acetic 1 2314.14 2314.14 H2S 1 2147.02 2147.02 Lactic 1 533.32 533.32 Residuals 26 2668.41 102.63

Response: taste Df Sum Sq Mean Sq Acetic 1 2314.14 2314.14 Lactic 1 1672.68 1672.68 H2S 1 1007.66 1007.66 Residuals 26 2668.41 102.63

