

A Statistics and Matrix Algebra Rosetta Stone

Overview

The knowledge that you already have in statistics can serve as an anchor for understanding concepts in matrix algebra and machine learning. Informative relationships can be found between covariance and inner products, beta and projection, and correlation and cosine distance. This is valuable, low-hanging fruit.

Below, we review basic statistics concepts from a more advanced perspective. Then we explore how those concepts can be represented with matrix algebra.

Statistics Review

Covariance

Covariance measures how much two variables co-vary. This is achieved by measuring the average of how much two variables deviate from their own means at the same time.

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

You will truly understand this when you realize that sorting two data columns, each from largest to smallest, maximizes their covariance. Furthermore, covariance, and therefore correlation, is a function of order; and we can shuffle data in order to impute a target level of correlation. That is the machinery of a copula.

R-Square and Correlation

The measure of covariance is in an awkward squared space. To make the measures of covariance interpretable, we can divide them by the total variance of X and Y.

$$R^2 = \frac{\text{Cov}(X, Y)^2}{\text{Var}(X) \cdot \text{Var}(Y)}$$

The resulting proportion of the covariation of X and Y as a fraction of all the variation of X and Y is the R-square. We can square root this value to move back to our one-dimensional metrics, and the result is correlation.

$$\text{Corr}(X, Y) = \sqrt{R^2}$$

We can summarize by saying that correlation is normalized covariance.

Variance and Standard Deviation

We can think of variance as the covariance of a variable with itself.

$$\text{Var}(X) = \text{Cov}(X, X)$$

Variance measures the average squared deviation from the mean:

$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

We can also think of variance as the original metric that measures dispersion as the average squared deviation. A consequence of that choice is that the arithmetic of risk will occur in squared space. For example, if you want to add measures of dispersion, you will add variances; and if you wanted to observe a fraction of explained dispersion, you will use ratios of variances. The act of square rooting these variances to calculate a standard deviation, moves us back to a one-dimensional metric and communicates the result.

$$\text{SD}(X) = \sqrt{\text{Var}(X)}$$

Beta

Beta can be calculated by identifying the shared covariation between X and Y as a fraction of all the variation in X:

$$\beta = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

Alternatively, Beta can be calculated starting with correlation, scaled by the ratio of standard deviations. Those standard deviations can be thought of as the relative step size for each of X and Y.

$$\beta = \text{Corr}(X, Y) \cdot \frac{\text{SD}(Y)}{\text{SD}(X)}$$

Matrix Algebra

Centering Data

Before proceeding, the trick that unlocks the following dual interpretations is work withing centered data. That is, the means have been subtracted from the data columns. From here onward, we will refer to data columns as vectors and it is assumed that the vectors represent centered data columns.

$$X = X^* - \bar{X}$$

Covariance, Correlation, and Beta

Covariance is proportional to the inner product of two vectors:

$$\text{Cov}(X, Y) = \frac{1}{n} \langle X, Y \rangle$$

The covariance matrix Σ for a multivariate dataset can be expressed as:

$$\Sigma = \frac{1}{n} X^T X$$

Correlation and beta are derived as follows:

$$\text{Corr}(X, Y) = \frac{\langle X, Y \rangle}{|X||Y|}, \quad \beta = \frac{\langle X, Y \rangle}{\langle X, X \rangle}$$

The norm of a vector corresponds to the standard deviation:

$$|X| = \sqrt{\frac{1}{n} \langle X, X \rangle}$$

If the mean has been estimated from data, then a degree of freedom has been lost, and the sum of deviations should be unitized by dividing by $n-1$. See the appendix ‘Where does that degree of freedom go?’ for detail.

Projection

The projection of a vector v onto w is:

$$p = sw, \quad s = \frac{\langle v, w \rangle}{\langle w, w \rangle}$$

Cosine similarity between two vectors is:

$$\text{Cosine Similarity} = \frac{\langle X, Y \rangle}{|X||Y|}$$

This is equivalent to correlation for centered data vectors.

Projection onto a Line

Let’s describe projection onto a line. We will do this twice, once in standard terminology and then with more applied terms.

Let's project vector \mathbf{v} onto vector \mathbf{w} . We are looking for a projection vector \mathbf{p} , as the location in \mathbf{w} that is closest to \mathbf{v} . \mathbf{p} is a scalar s multiple of \mathbf{w} .

$$\mathbf{p} = s\mathbf{w}$$

Let \mathbf{r} represent the residual distance from \mathbf{v} to \mathbf{p} . Visually, \mathbf{r} drops down from \mathbf{v} onto \mathbf{p} .

Mathematically, $\mathbf{v} + \mathbf{r} = \mathbf{p}$ and we choose \mathbf{p} so that \mathbf{r} falls onto it orthogonally, i.e., $\mathbf{r} \cdot \mathbf{p} = 0$.

Derivation

$$\mathbf{p} = s\mathbf{w}$$

$$\mathbf{v} + \mathbf{r} = \mathbf{p}$$

$$\mathbf{r} \cdot \mathbf{w} = 0$$

We solve for the scalar (s):

$$\mathbf{v} + \mathbf{r} = s\mathbf{w}$$

Multiply everything by \mathbf{w} :

$$\mathbf{v} \cdot \mathbf{w} + \mathbf{r} \cdot \mathbf{w} = s(\mathbf{w} \cdot \mathbf{w})$$

Since $\mathbf{r} \cdot \mathbf{w} = 0$, this simplifies to:

$$s = \frac{\mathbf{v} \cdot \mathbf{w}}{\mathbf{w} \cdot \mathbf{w}}$$

Generalization to Subspaces

\mathbf{w} does not need to be a vector; it could be a multi-dimensional subspace W .

Beta in Context

Starting with a data vector \mathbf{y} and predictor vector \mathbf{x} , we want to project \mathbf{y} onto $\beta\mathbf{x}$. Notice the difference; we are projecting \mathbf{y} onto $\hat{\mathbf{y}}$ i.e. $\beta\mathbf{x}$. The math proceeds similarly after defining \mathbf{r} as the residual vector:

$$\mathbf{y} + \mathbf{r} = \beta\mathbf{x}$$

$$\mathbf{x} \cdot \mathbf{y} + \mathbf{x} \cdot \mathbf{r} = \beta(\mathbf{x} \cdot \mathbf{x})$$

Since $\mathbf{x} \cdot \mathbf{r} = 0$:

$$\mathbf{x} \cdot \mathbf{y} = \beta(\mathbf{x} \cdot \mathbf{x})$$

$$\beta = \frac{\mathbf{x} \cdot \mathbf{y}}{\mathbf{x} \cdot \mathbf{x}}$$

This schema generalizes to multiple linear regression:

$$\begin{aligned} \mathbf{y} + \mathbf{r} &= \mathbf{X}\beta \\ \mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{r} &= \mathbf{X}^T \mathbf{X} \beta \\ \text{Since } \mathbf{X}^T \mathbf{r} &= 0: \\ \mathbf{X}^T \mathbf{y} &= \mathbf{X}^T \mathbf{X} \beta \\ \beta &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

This time, we projected into the subspace of \mathbf{X} . We identified the subspace of \mathbf{X} that is closest to \mathbf{y} . Following the principle that the shortest distance between two points is a straight line, we have identified the closest subspace of \mathbf{X} by projecting onto it. This reasoning is what allows for this closed-form solution to Ordinary Least Squares.

Cosine Similarity and Correlation

Cosine similarity is used in natural language processing (NLP), recommendation systems, and data clustering. In NLP, it measures the similarity between texts, such as comparing word embeddings, sentences, or documents, for tasks like semantic search and text classification. In recommendation systems, cosine similarity helps identify users or items with similar preferences. It is also used in clustering algorithms, especially in high-dimensional and sparse datasets like social networks or document corpora.

Cosine is defined as:

$$\cos(\theta) = \frac{\text{adjacent}}{\text{hypotenuse}}$$

Replacing terminology with \mathbf{x} and \mathbf{y} :

$$\cos(\theta) = \frac{|\mathbf{x}|}{|\mathbf{y}|}$$

Using the projection formula, $s = \frac{\mathbf{x} \cdot \mathbf{y}}{\mathbf{x} \cdot \mathbf{x}}$, we derive:

$$\cos(\theta) = \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$$

This shows how cosine similarity is equivalent to correlation for centered data.

Projection and Spatial Algorithms

Projection forms the basis for algorithms like Linear Discriminant Analysis (LDA) and Support Vector Machines (SVM).

LDA Objective Function

$$\frac{(\text{mean projection for A} - \text{mean projection for B})^2}{\text{variance projection for A} + \text{variance projection for B}}$$

SVM

Identify \mathbf{w} such that a hyperplane splits the categories.

Understanding these principles deepens comprehension of regression, LDA, and SVM.

Appendix: Where Does That Degree of Freedom Go?

If you know the population mean μ , then the formula for variance is:

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{n}$$

This formula works because μ is a known constant. However, if \bar{x} (the sample mean) is used instead of μ , there is a problem.

Specifically, \bar{x} comes from the data and “uses up” some of the information. In fact, it uses $1/n$ of the data per data point. To demonstrate, consider this estimation of variance for three data points:

$$(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + (x_3 - \bar{x})^2$$

Do you see anything strange? Maybe not yet, but let’s expand \bar{x} :

$$\bar{x} = \frac{x_1 + x_2 + x_3}{3}$$

Substituting this back, the expression becomes:

$$\left(x_1 - \frac{x_1 + x_2 + x_3}{3}\right)^2 + \left(x_2 - \frac{x_1 + x_2 + x_3}{3}\right)^2 + \left(x_3 - \frac{x_1 + x_2 + x_3}{3}\right)^2$$

Notice that for each of these measures of squared deviation, the data point x_i is going to subtract $1/3$ of itself and get a perfect 0 for part of the calculation. This 0 was supposed to be a measure of dispersion, and it would have been if we had been able to use μ .

The consequence of using \bar{x} is that we lose $1/n$ of a measure of dispersion for each data point, and this happens n times. Ultimately, we lose one full unit of dispersion, also known as one degree of freedom.

This is why the formula for sample variance includes an adjustment:

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$

Here, $n - 1$ accounts for the lost degree of freedom due to the estimation of \bar{x} from the data.