

1 Assumptions

We treat the data as if it were a point cloud, meaning that at each time step, t , the values of all channels at time t as a vector are a point. We treat all points as if they are independent and are drawn from the same distribution.

2 Problem Setting

Say that we have n channels, c_1, \dots, c_n , and that $C_i(t)$ is the value of all channels but i at a given time t .

For each channel c_i we have a map $f_i : \mathbb{R}^{n-1} \rightarrow \mathbb{R}$ that, given the state of other channels at a given point in time, predicts the value of channel i . For each f_i , we also have the metric ρ_i , where

$$\rho_i = \frac{\text{Cov}_t(f_i(C_i(t)), c_i(t))}{\sqrt{\text{Var}_t(f_i(C_i(t))) \text{Var}_t(c_i(t))}},$$

which is the Pearson correlation coefficient between the predicted channel, given by $f_i(C_i(t))$ for a single point in time, and the actual channel c_i . Note that $|\rho_i|$ may not be large for some channels, meaning that our function f_i does a poor job fitting this channel. However, our model in detecting anomalies will account for the fact that many fits may be imperfect.

Then, we will use our model and some additional statistics to check for anomalies in real-time. At a high level, our model will calculate the correlation between how our model predicts what the sample will be versus what the sample actually is. For a window of length k , meaning entries from $t_0 - k + 1, \dots, t_0$, where t_0 is the current time, we will find r_i , the test statistic, given by

$$r_i = \frac{\text{Cov}_{t_0-k < t \leq t_0}(f_i(C_i(t)), c_i(t))}{\sqrt{\text{Var}_{t_0-k < t \leq t_0}(f_i(C_i(t))) \text{Var}_{t_0-k < t \leq t_0}(c_i(t))}},$$

We will show how to calculate this test statistic efficiently (in $\Theta(1)$ time amortized per window) and will show how to reject a sample as an anomaly using probabilistic methods.

For the rest of the discussion, we will fix some t_0 and some k . For notation purposes, let $X_j = f_i(C_i(t_0 - k + j))$ and $Y_j = c_i(t_0 - k + j)$ for $1 \leq j \leq k$, meaning that X_j and Y_j are the j th entries of the predicted and actual samples of the time series. Then, let A and B be the normalized version of X and Y , meaning that

$$A_i = \frac{X_i - \bar{X}}{\sigma_X}, \quad B_i = \frac{Y_i - \bar{Y}}{\sigma_Y},$$

such that

$$\mathbb{E}[A] = \mathbb{E}[B] = 0, \quad \text{Var}(A) = \text{Var}(B) = 1.$$

It is important to note that this normalization is really just a “trick” to simplify the calculations. We additionally note that it does not matter whether we use the sample variance (multiplied by $\frac{k}{k-1}$ or not as long as we are consistent and also use the sample covariance, as both result multiplying the top and bottom of the expression for r_i by the same value. Now, with A and B , we have that the expression for r_i is

$$r_i = \frac{\text{Cov}(A, B)}{\text{Var}(A) \text{Var}(B)},$$

but since the $\text{Var}(A) = \text{Var}(B) = 1$, this is just

$$r_i = \text{Cov}(A, B).$$

Then, by the definition of the covariance and since we know the values of A and B this is just

$$r_i = \text{Cov}(A, B) = \mathbb{E}[AB] = \frac{1}{k} \sum_{i=1}^k A_i B_i,$$

which, interestingly, is the cosine-distance if we treat these as scalars (an aside: since $\text{Var}[A] = \text{Var}[B] = 1$, this says $\frac{1}{k} \|A\|_2^2 = \frac{1}{k} \|B\|_2^2 = 1$, meaning that $\|A\|_2 = \|B\|_2 = \sqrt{k}$, so the cosine distance $\frac{A \cdot B}{\|A\|_2 \|B\|_2}$ is $\sum_{i=1}^k A_i B_i / \sqrt{k * k} = \frac{1}{k} \sum_{i=1}^k A_i B_i = r_i$).

Now, we want to run a statistical test to determine if the window ending at t_0 represents an anomaly. Then, we use the two hypotheses:

$$H_0 : \mu_{r_i} = \rho_i, \quad H_a : \mu_{r_i} \neq \rho_i.$$

The null hypothesis, H_0 , represents that the correlation between this window of the data is the same as the correlation that we would expect. This then represents that the data is as expected and that there is no anomaly. H_a represents some sort of anomaly in the window. Our goal will then be to calculate the probability p that this window has correlation r_i given that H_0 is true, and if p is “low enough” (to be discussed later) then we can say with confidence that there is an anomaly.

Then, to calculate this probability, we want to know the distribution of r_i . We can think of A_i and B_i themselves as being random variables, so $A_i B_i$ is a random variable, and by our assumptions these are independent and identically distributed. Then, since r_i is a sum of independent indentially distributed random variables, the Central Limit Theorem applies, which says that in the limit (as $k \rightarrow \infty$), that r_i is normally distributed. This gives us a good approximation for the distribution of r_i . Then, since we assume that $\mu = \rho$, the only parameter we need estimate is the standard deviation of r_i . We have that this is

$$\text{Var}(r_i) = \text{Var}\left(\frac{1}{k} \sum_{i=1}^k A_i B_i\right) = \frac{1}{k^2} \sum_{i=1}^k \text{Var}(A_i B_i),$$

by linearity, but since $A_i B_i$ are all identically distributed, their variances are the same, so we can write this as

$$= \frac{1}{k} \text{Var}(AB) = \mathbb{E}((AB - r_i)^2) = \frac{1}{k} \left(\mathbb{E}((AB)^2) - r_i^2 \right),$$

by a well-know formula for variance. To calculate $\mathbb{E}((AB)^2)$, we just do

$$\mathbb{E}((AB)^2) = \sum_{i=1}^k \frac{1}{k} (A_i B_i)^2.$$

Additionally, to get the sample variance, we multiply by $\frac{k}{k-1}$, so the full formula for $S_{r_i}^2$ is

$$\text{Var}(r_i) = \frac{1}{k-1} \left(\sum_{i=1}^k (A_i B_i)^2 - r_i^2 \right).$$

Then, this means that

$$\frac{r_i - \rho}{\sqrt{S_{r_i}^2}}$$

is distributed approximately normally, which means that given some other standard normal random variable N that the probability of drawing r_i randomly under the assumption that $\mu_{r_i} = \rho$ is

$$p \approx P \left(|N| \geq \left| \frac{r_i - \rho}{\sqrt{S_{r_i}^2}} \right| \right) = 2P \left(N \leq - \left| \frac{r_i - \rho}{\sqrt{S_{r_i}^2}} \right| \right),$$

which we can calculate using existing normal CDF functions.

3 Finding r_i efficiently

We can see from above that we can calculate r_i in $\Theta(k)$ time, where k is the size of the window, which isn't bad. However, if we want the anomaly detection system to work in real-time, we want to calculate r_i faster than this. We will show a method for calculating r_i in amortized $\Theta(1)$ time and using $\Theta(k)$ memory, assuming that we are calculating r_i for every window of size k .

The method works by using two different queues that store the last k values and additionally by maintaining a number of different accumulators that store some value that changes as we move from left to right.

For calculating r_i , we have

$$\begin{aligned}
 r_i &= \text{Cov}(A, B) = \text{Cov}\left(\frac{X_i - \bar{X}}{\sigma_X}, \frac{Y_i - \bar{Y}}{\sigma_Y}\right) \\
 &= \frac{1}{\sigma_X \sigma_Y} \text{Cov}(X_i - \bar{X}, Y_i - \bar{Y})
 \end{aligned}$$