# Correlations:

## The Tango of Time Series

Lance Hester
Learning Tuesday
06/02/2020

# What you Should Take Away  Today
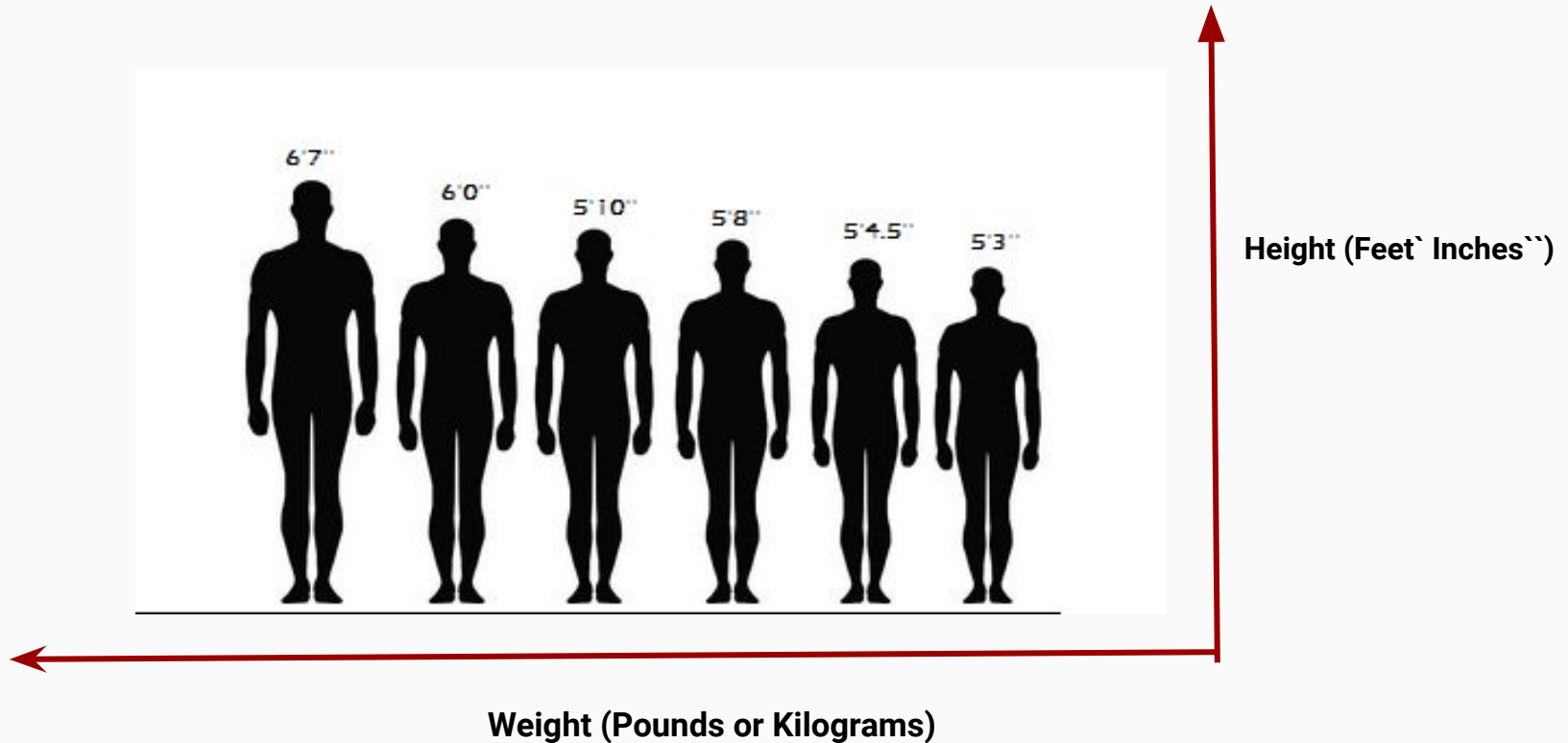
Getting to know your time series by:

- Comparing it to itself (auto = "self")
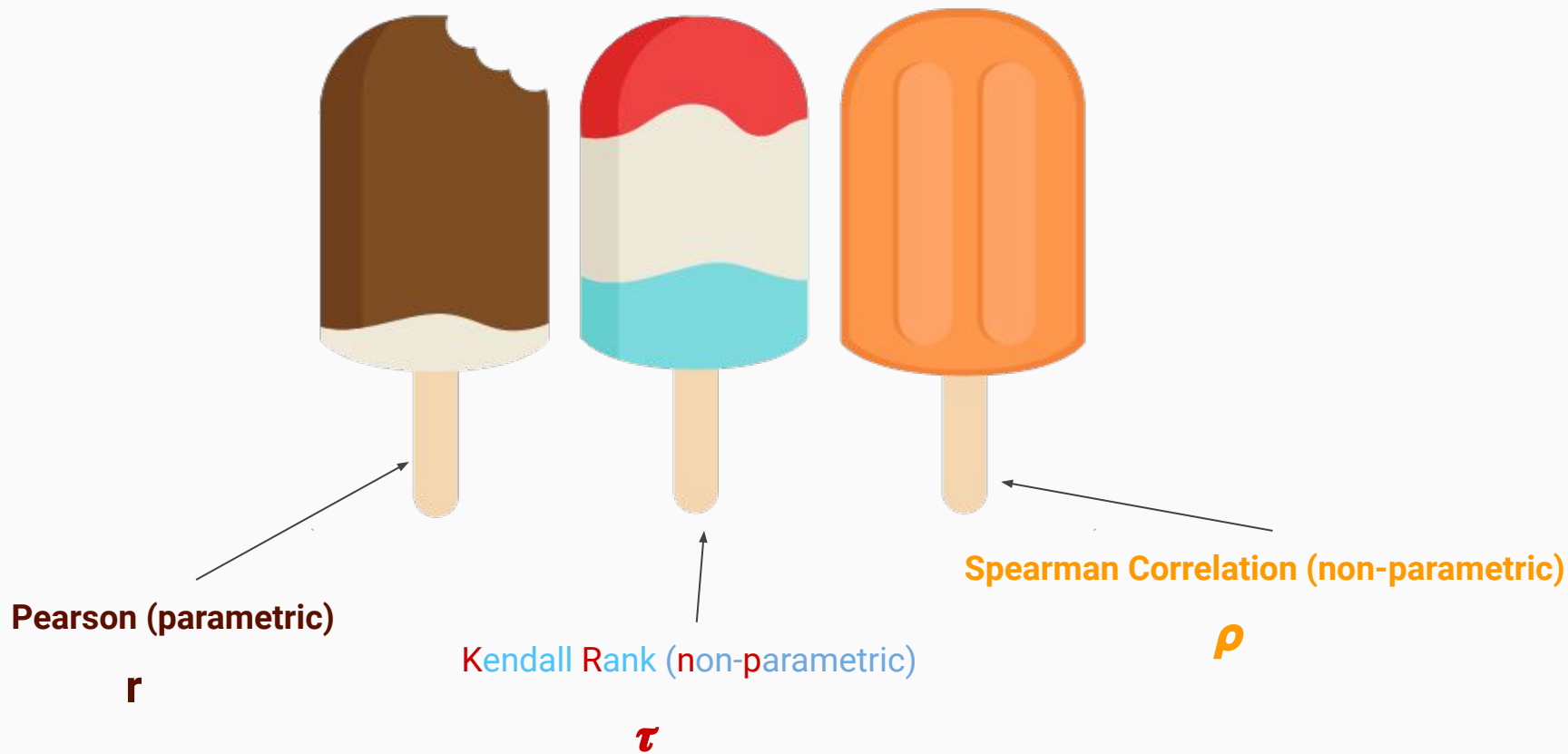- Comparing other time series (cross)

# Correlation

**Definition** *cor·re·la ·tion*

The degree of correspondence or relationship between two variables.

# Correlation of Two Metrics/Time Series

# Formula Flavors



Pearson (parametric)

$r$

Kendall Rank (non-parametric)

$\tau$

Spearman Correlation (non-parametric)

$\rho$

# Pearson Correlation Formula

$$\rho_{X,Y} = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\text{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$
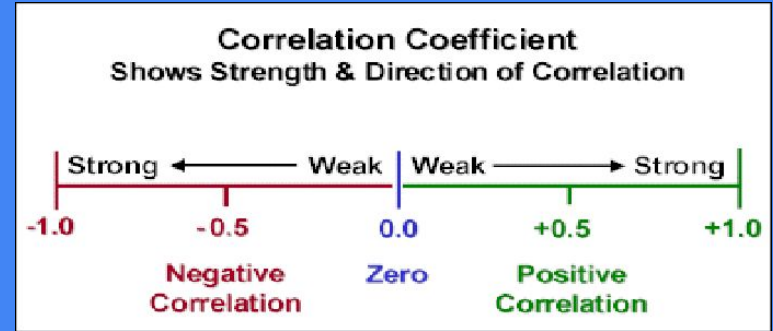
$$\rho_{X,Y} = \frac{\text{E}(XY) - \text{E}(X)\,\text{E}(Y)}{\sqrt{\text{E}(X^2) - \text{E}(X)^2} \cdot \sqrt{\text{E}(Y^2) - \text{E}(Y)^2}}$$

# Correlation Coefficient r

(quantitative measure)
(aka, a number!)

**Correlation Coefficient**
**Shows Strength & Direction of Correlation**

| Strong ← Weak | Weak → Strong |
|---|---|
| -1.0   -0.5   0.0 | +0.5   +1.0 |

Negative Correlation    Zero    Positive Correlation

## Strength

## Direction

# Thanks!

Provided good examples of Correlations which I Julia-ized and added some additional information.

https://anomaly.io/detect-anomalies-in-correlated-time-series/index.html

https://anomaly.io/understand-auto-cross-correlation-normalized-shift/index.html#/cross_correlation

https://anomaly.io/detect-correlation-time-series/index.html

# Cross-Correlation:

# It Takes Two!

```
a = [1,2,-2,4,2,3,1,0];
b = [2,3,-2,3,2,4,1,-1];
c = [-2,0,4,0,1,1,0,-2];
```

$$corr(x, y) = \sum_{n=0}^{n-1} x[n] * y[n]$$

## TS plots

$corr(a, b) = 1 * 2 + 2 * 3 + -2 * -2 + 4 * 3 + 2 * 2 + 3 * 4 + 1 * 1 + 0 * -1$
$= 41$

$corr(a, c) = 1 * -2 + 2 * 0 + -2 * 4 + 4 * 0 + 2 * 1 + 3 * 1 + 1 * 0 + 0 * -2$
$= -5$

# Issues with Raw Cross-Correlation Calculations

1. Can't really grasp value of cross_ab vs cross_ac significance
2. Want a&b or a&c to have similar amplitudes - might misread correlations.

$$corr(a, a/2) = 1 * (1/2) + 2 * (2/2) + -2 * (-2/2) + 4 * (4/2) + 2 * (2/2), 3 * (3/2) + 1 * (1/2) + 0 * (0/2)$$
$$= 19.5$$

3. Have to ensure that std deviation values are finite and positive.

# "Solution: Normalize the Values"

Normalization alleviates these issues so we can compare.

$$norm\_corr(x, y) = \frac{\sum_{n=0}^{n-1} x[n] * y[n]}{\sqrt{\sum_{n=0}^{n-1} x[n]^2 * \sum_{n=0}^{n-1} y[n]^2}}$$

$$norm\_corr(x, y) = PearsonCorrelationCoefficient$$

Using this formula let's compute the normalized cross-correlation of ab and ac.

$$norm\_corr(a, b) = \frac{1 * 2 + 2 * 3 + -2 * -2 + 4 * 3 + 2 * 2 + 3 * 4 + 1 * 1 + 0 * -1}{\sqrt{(1 + 4 + 4 + 16 + 4 + 9 + 1 + 0) * (4 + 9 + 4 + 9 + 4 + 16 + 1 + 1)}}$$

$$= \frac{41}{\sqrt{(39) * (48)}}$$

$$= 0.947$$

$$norm\_corr(a, c) = \frac{1 * -2 + 2 * 0 + -2 * 4 + 4 * 0 + 2 * 1 + 3 * 1 + 1 * 0 + 0 * -2}{\sqrt{(1 + 4 + 4 + 16 + 4 + 9 + 1 + 0) * (4 + 0 + 16 + 0 + 1 + 1 + 0 + 4)}}$$

$$= \frac{-5}{\sqrt{(39) * (26)}}$$

$$= -0.157$$

```
norm_corr_ab = sum(a .* (b)) / sqrt(sum(a.^2) * sum(b.^2)); # equals  0.947
norm_corr_ac = sum(a .* (c)) / sqrt(sum(a.^2) * sum(c.^2)); # equals -0.157
```

# Quick Check to Show Normalization Works

```
In [4]:   1  # Normalized norm_corr(a,a) = 1:
          2  proof = sum(a .* (a)) / sqrt(sum(a.^2) * sum(a.^2))

Out[4]: 1.0
```

```
In [5]:   1  # Normalized norm_corr(a,-a) = -1:
          2  proof = sum(a .* (-a)) / sqrt(sum(a.^2) * sum(((-a).^2)))

Out[5]: -1.0
```

```
In [6]:   1  # Normalized cross-correlation can detect the correlation of two signals with different amplitudes:
          2  # Notice we have perfect correlation between signal A and the same signal with half the amplitude!
          3
          4  proof = sum(a .* (a/2)) / sqrt(sum(a.^2) * sum((a./2).^2))

Out[6]: 1.0
```

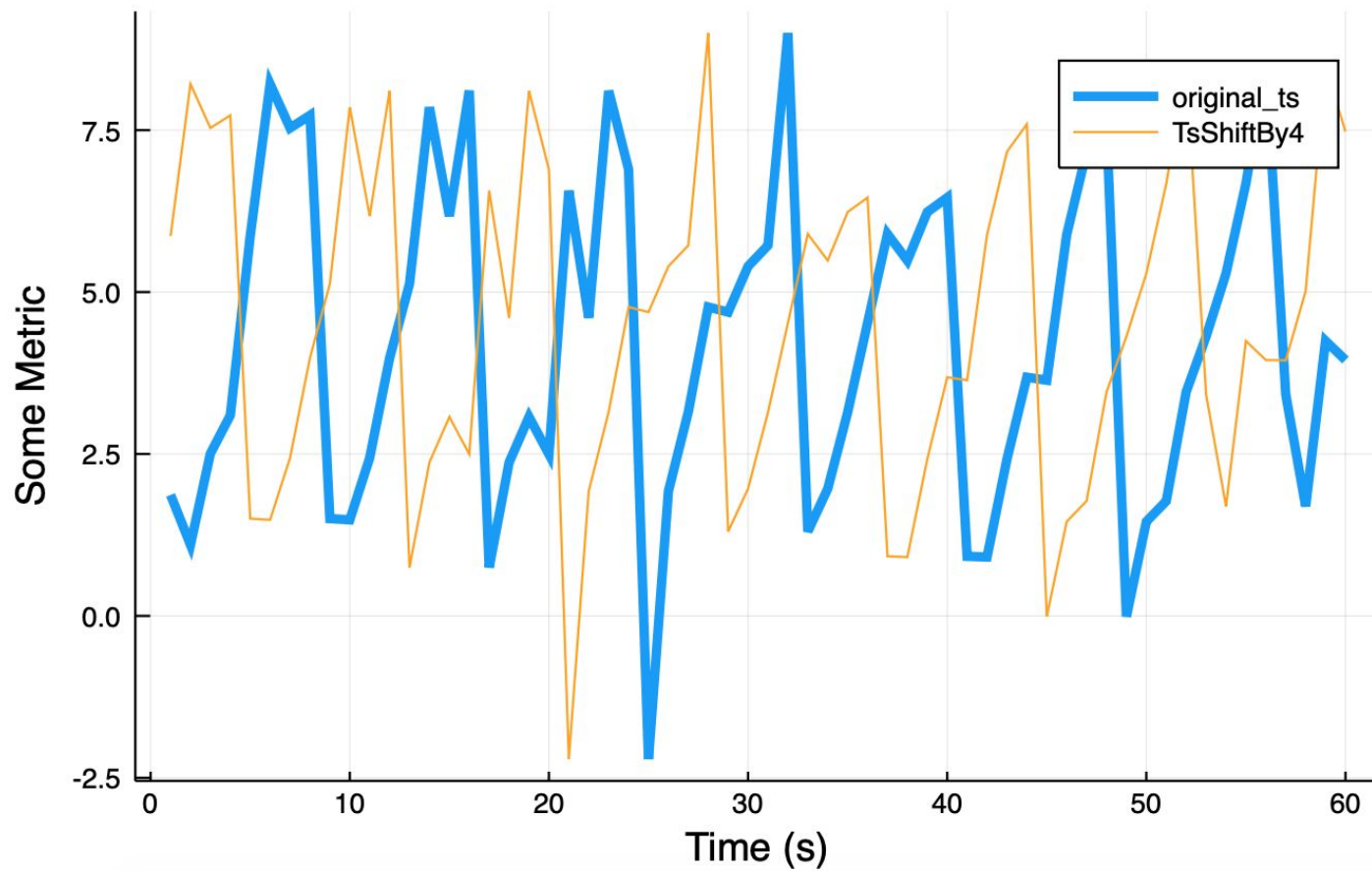Autocorrelation:

One is the Loneliest Number!

Original Signal
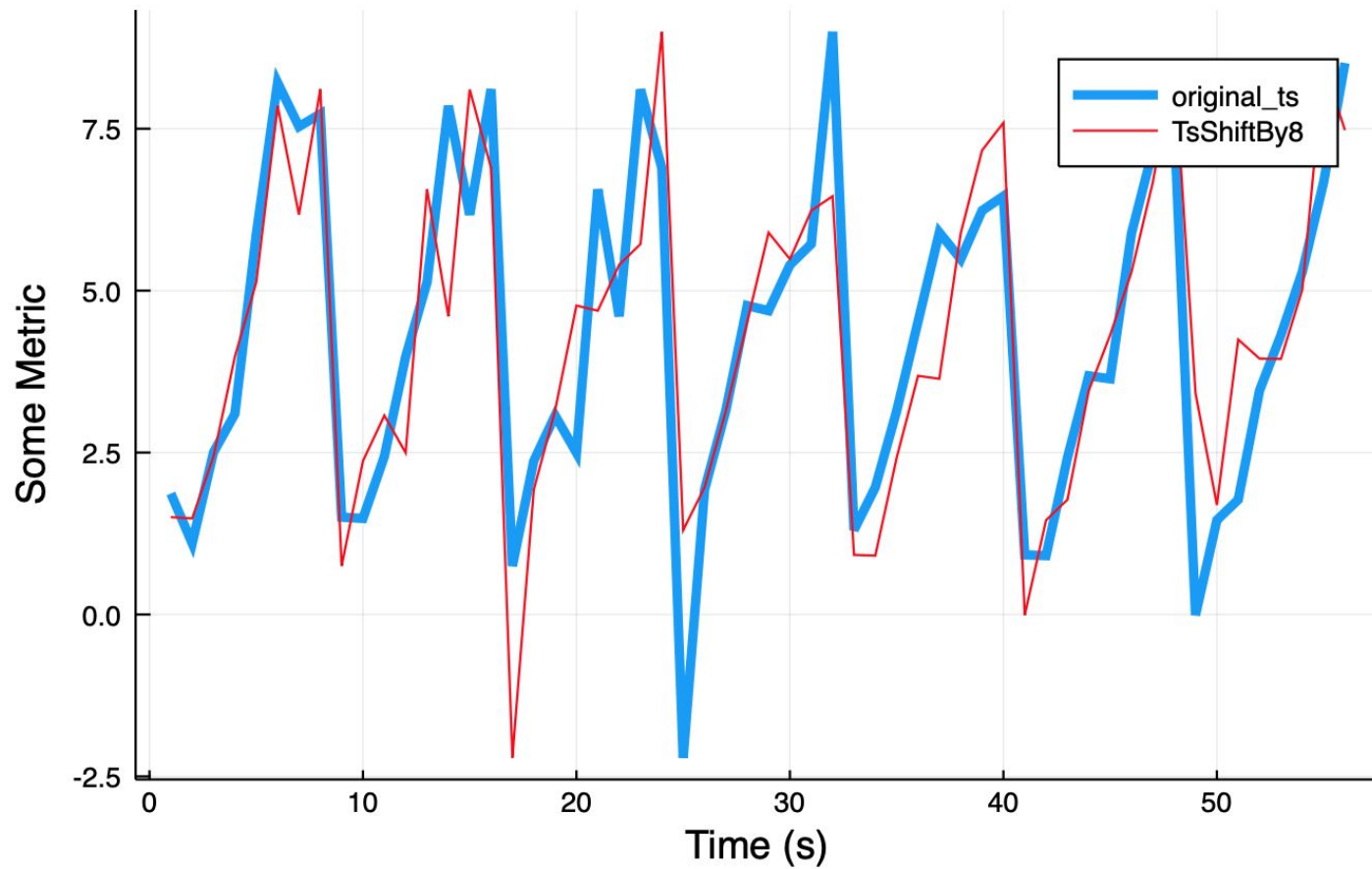
Original Signal vs TsShiftBy4

Original Signal vs TsShiftBy8

# Autocorrelation

1. Comparison of time series with itself at different times
2. Auto-Correlation detects repeating patterns or "**seasonality**" !
3. Auto-Correlation answers questions like:
   a. Can we see some weekly pattern?
   b. Is today similar to last week today?

# Unnormalized Autocorrelation

```
1   #Computing the Correlations -- here autocorrelations (i.e., multiplying and Summing the two
2
3   corr_shift4 = sum(autoSignalFour .* autoSignalShiftFour); # equals 948.4089186791925
4   corr_shift8 = sum(autoSignalEight .* autoSignalShiftEight); #equals 1336.0693024826921
```

Leads us to believe we have detected possible seasonality at 8

# Normalized Autocorrelation Makes it Obvious

```
1  norm_auto_shift4 = sum(autoSignalFour .* autoSignalShiftFour) / sqrt(sum(autoSignalFour.^2)
2  norm_auto_shift8 = sum(autoSignalEight .* autoSignalShiftEight) / sqrt(sum(autoSignalEight.
```

```
1  # norm_auto_shift4 = 0.6227933971623315
2  # norm_auto_shift8 = 0.9602671052926668     == Normalized autocorrelation makes it very obviou
```

```
In [12]:    1  rng = MersenneTwister(1234);
            2  numRepeats = 8;
            3  autoSignal = repeat(collect(1.0:8.0), numRepeats) + randn!(rng, zeros(numRepeats*8));
            4  x = collect(1:length(autoSignal))
            5  plot(x,autoSignal,
            6      title = "Original Signal",
            7      label =["original_ts"],
            8      xlabel="Time (s)",
            9      ylabel="Some Metric",
           10      legend=:bottomright,
           11      linecolor = 1,
           12      lw =2)
```
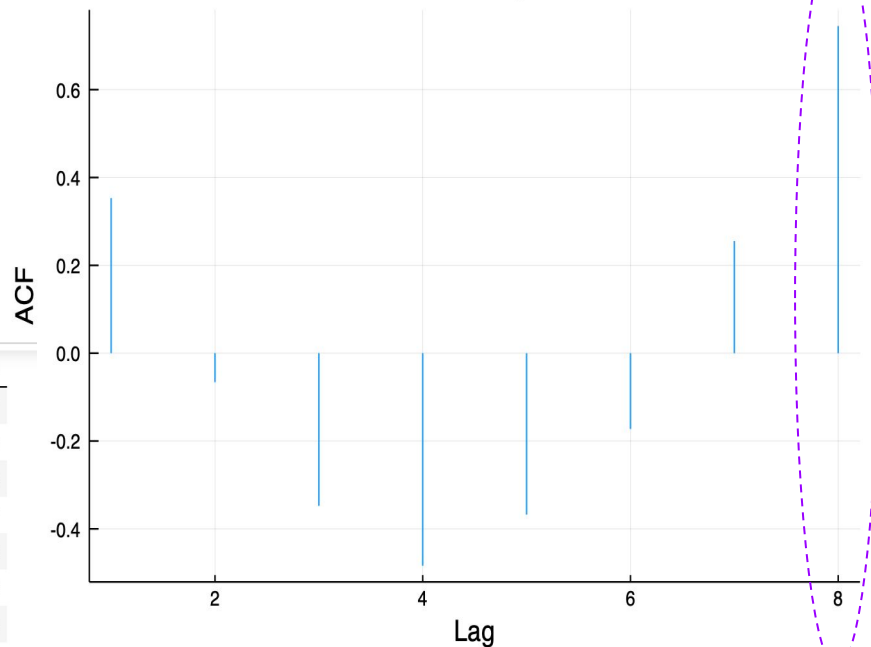
```
In [13]:    1  auto_r = autocor(autoSignal, collect(1:8));
```

```
l5]:    1  plot(acf_df[!,:lag], acf_df[!,:acfValue], line = :sticks, legend=false,
        2      xlabel="Lag", ylabel="ACF", title = "ACF of Signal", )
```



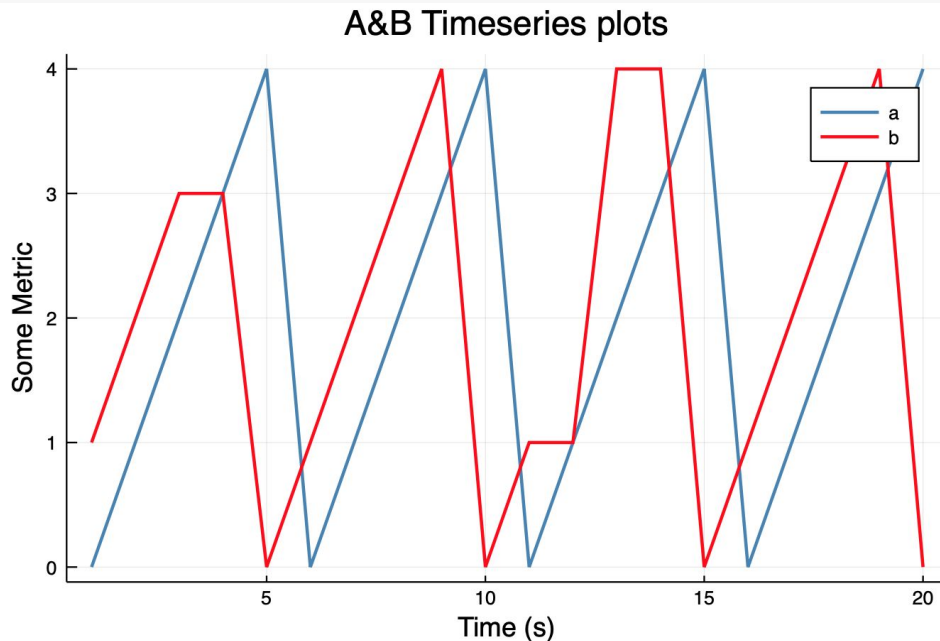| | Int64 | Float64 |
|---|---|---|
| 1 | 1 | 0.353117 |
| 2 | 2 | -0.0659904 |
| 3 | 3 | -0.347516 |
| 4 | 4 | -0.48397 |
| 5 | 5 | -0.36758 |
| 6 | 6 | -0.1726 |
| 7 | 7 | 0.255374 |
| 8 | 8 | 0.744855 |

# One More Bite at Time-Shifting:

# Back to Cross-Correlations

# Cross-Correlation with Time Shifts

1. Check to see if one signal compared to another
   a. Lags (delays) - move elements to the right (t-lag)
   b. Leads (advancing) - move elements to left (t+lead)


2. Find the best time-shift at which time signals are correlated
   a. In autocorrelation, that is lag=lead=0 (most energy at perfect overlap)

```
In [16]:    1  a = [0,1,2,3,4,0,1,2,3,4,0,1,2,3,4,0,1,2,3,4]
            2  b = [1,2,3,3,0,1,2,3,4,0,1,1,4,4,0,1,2,3,4,0]
            3  time = collect(1:length(a));
            4
            5  plot(time, a, label="a", lw=2, linecolor=:steelblue, xlabel="Time (s)",
            6       ylabel="Some Metric", title = "A&B Timeseries plots")
            7  plot!(time, b, label="b", lw=2, linecolor=:red)
            8
            9
```
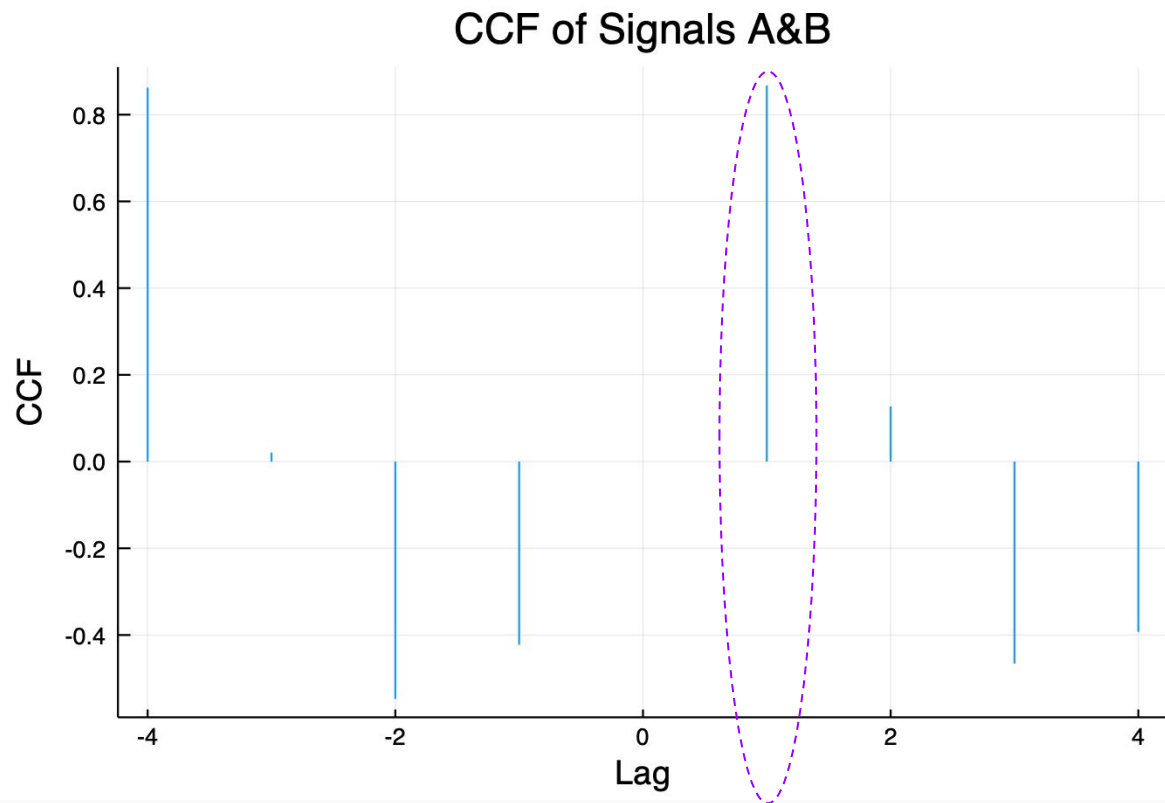


A&B Timeseries plots

```
In [72]:  1
          2  r = crosscor(b, a, collect(-4:4))
```

Out[21]:

| | ccf_lag | ccfValue |
|---|---|---|
| | Int64 | Float64 |
| 1 | -4 | 0.86232 |
| 2 | -3 | 0.0210021 |
| 3 | -2 | -0.547289 |
| 4 | -1 | -0.422512 |
| 5 | 0 | -3.29181e-17 |
| 6 | 1 | 0.867262 |
| 7 | 2 | 0.127248 |
| 8 | 3 | -0.465752 |
| 9 | 4 | -0.392862 |



CCF of Signals A&B

Lag = 1 => Best Correlation

Okay Smarty,
How Does Correlation
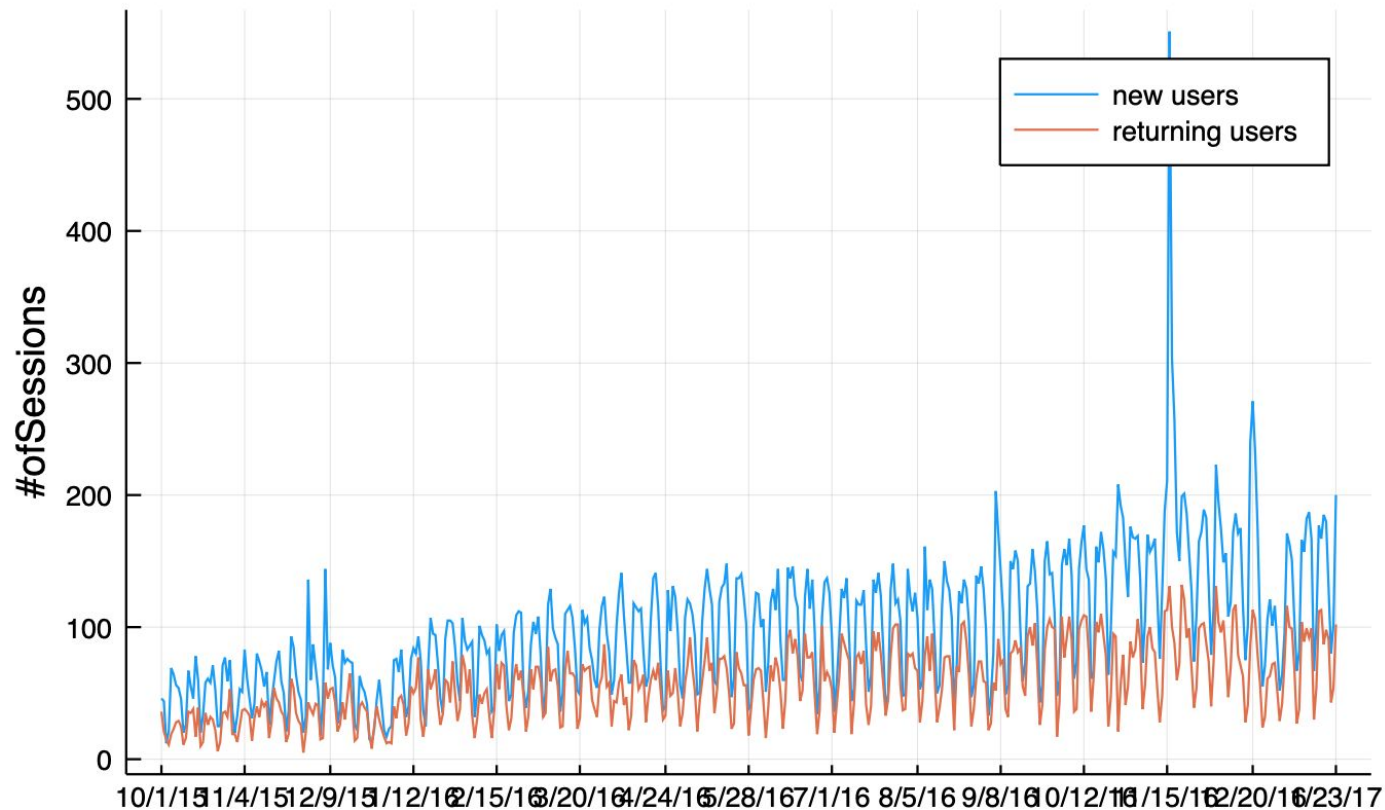Help with Anomaly
Detection?

# Checkout: Key Performance Indicators (KPIs)

New Users vs Returning Users

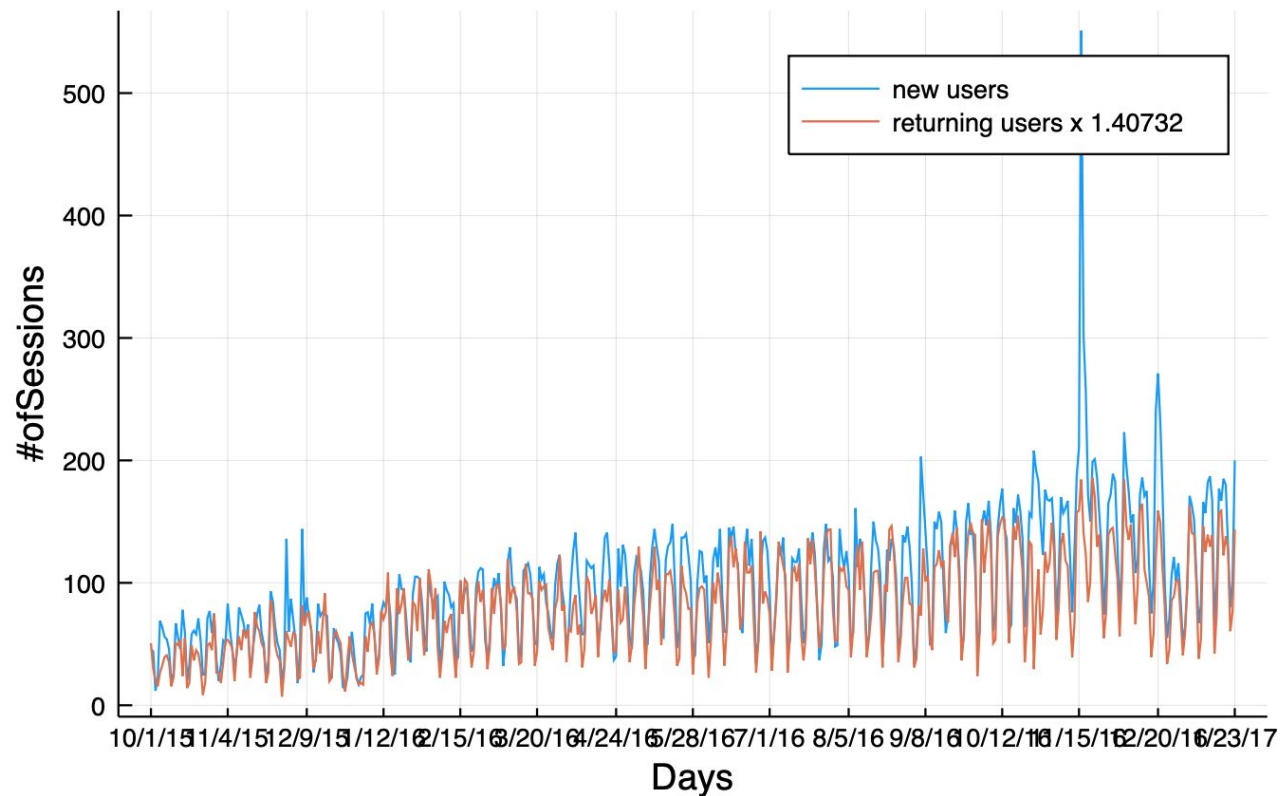to a website

New Users vs Returning Users

```
In [27]:    1   crosscor(new_data_df[!,:Sessions], return_data_df[!,:Sessions], [0])

Out[27]:    1-element Array{Float64,1}:
            0.8371268696198125
```

Show high correlation > 0.7

# New Users vs Returning Users x 1.40732
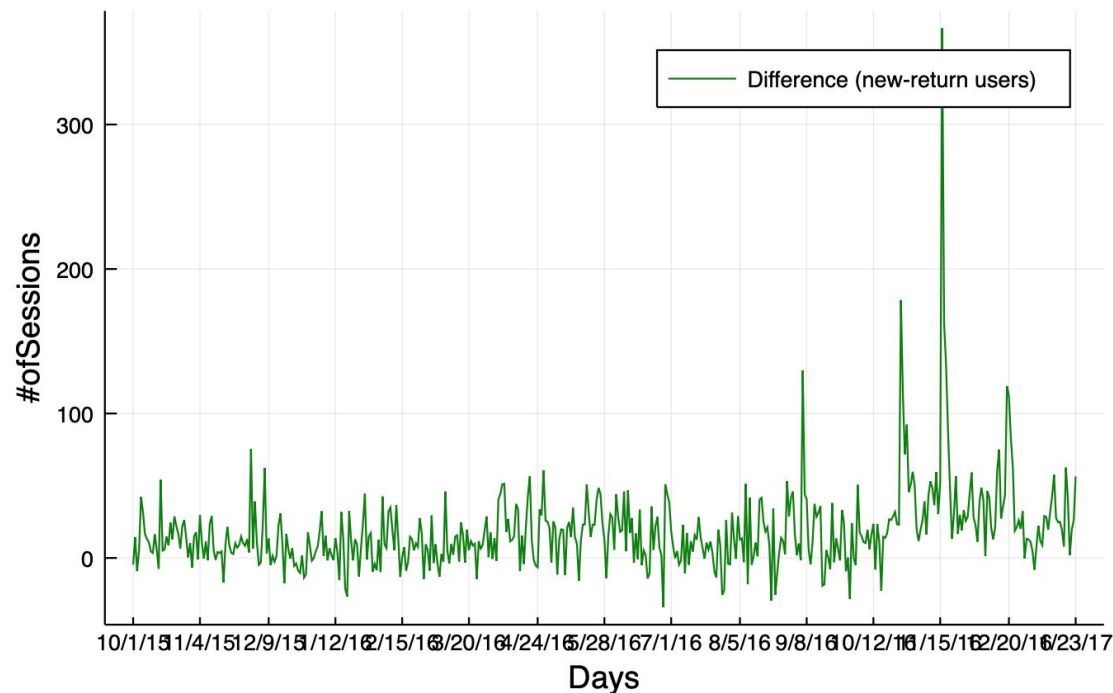


```
In [28]:  1  # Subtract the time Series
          2  multi = sum(new_data_df[!,:Sessions]/return_data_df[!,:Sessions])/ length(return_data_df[!,:Sessions])
          3  multi #display 1.40732
          4  align_return_data = return_data_df[!,:Sessions].*multi;
```

```
In [31]:   1  subtractTs = new_data_df[!,:Sessions] - align_return_data;
           2
```

```
In [32]:   1  plot(return_data_df[!,:Day_Index], subtractTs,
           2      xlabel="Days", ylabel="#ofSessions", title = "Diff New Users and Returning Users",
           3      label="Difference (new-return users)", linecolor=:green)
```
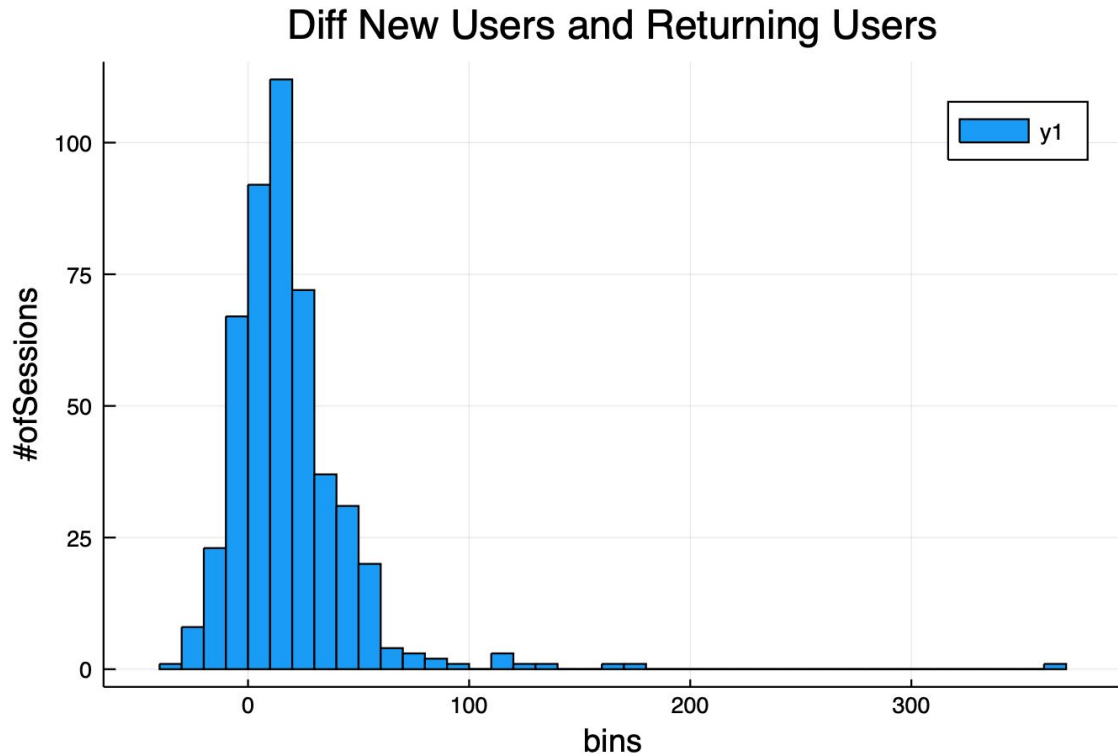
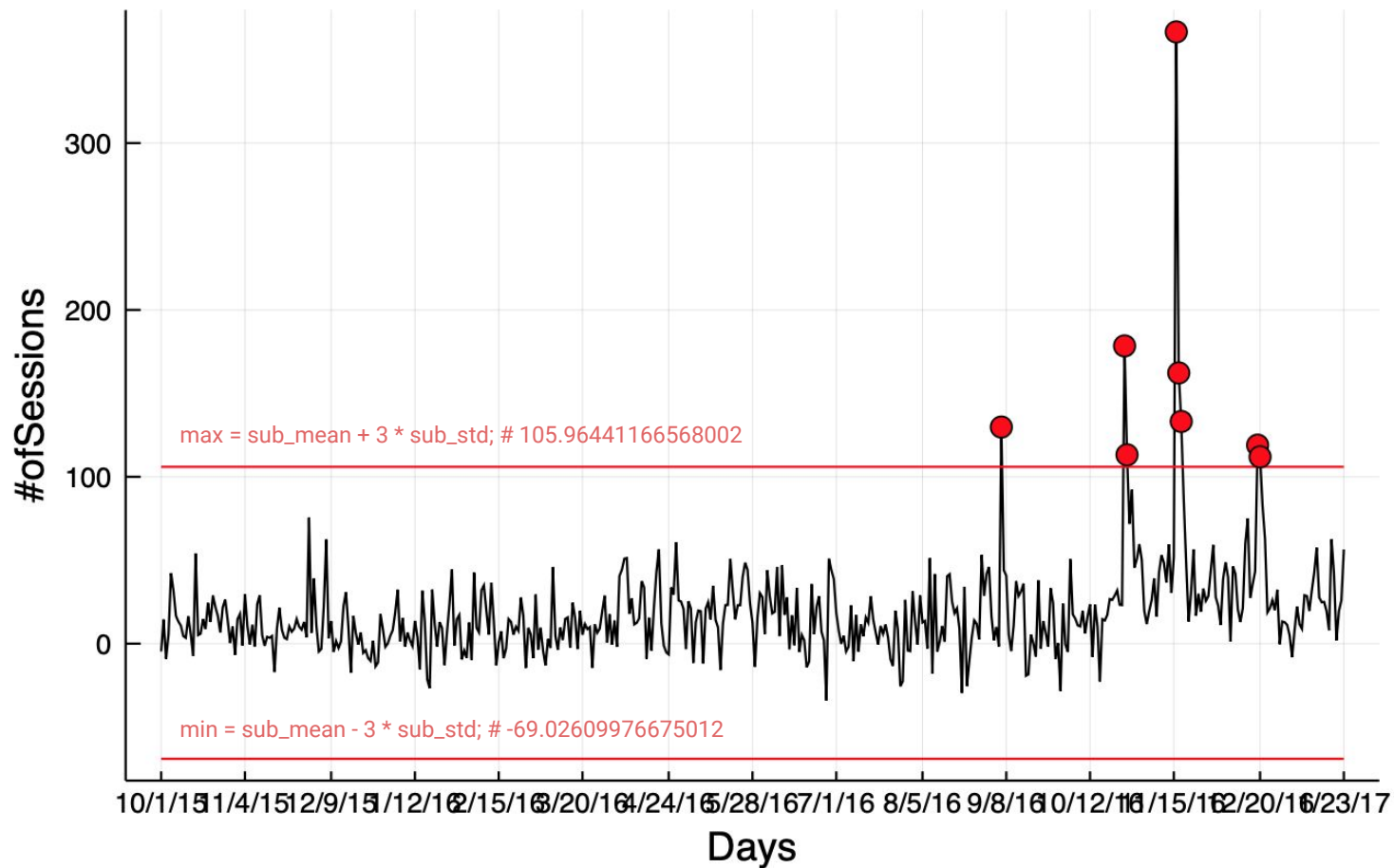Out[32]:



Diff New Users and Returning Users

```
1  # Finding Outliers in Correlated Time Series
2  histogram(subtractTs, xlabel="bins", ylabel="#ofSessions", title = "Diff New Users and Returning Users
```

## Diff New Users and Returning Users

3-Sigma Outliers

max = sub_mean + 3 * sub_std; # 105.96441166568002

min = sub_mean - 3 * sub_std; # -69.02609976675012

| | index | value | date |
| --- | --- | --- | --- |
| | Int64 | Float64 | String |
| 1 | 342 | 129.819 | 9/6/16 |
| 2 | 392 | 178.446 | 10/26/16 |
| 3 | 393 | 113.19 | 10/27/16 |
| 4 | 413 | 366.641 | 11/16/16 |
| 5 | 414 | 162.268 | 11/17/16 |
| 6 | 415 | 133.155 | 11/18/16 |
| 7 | 446 | 118.97 | 12/19/16 |
| 8 | 447 | 111.972 | 12/20/16 |