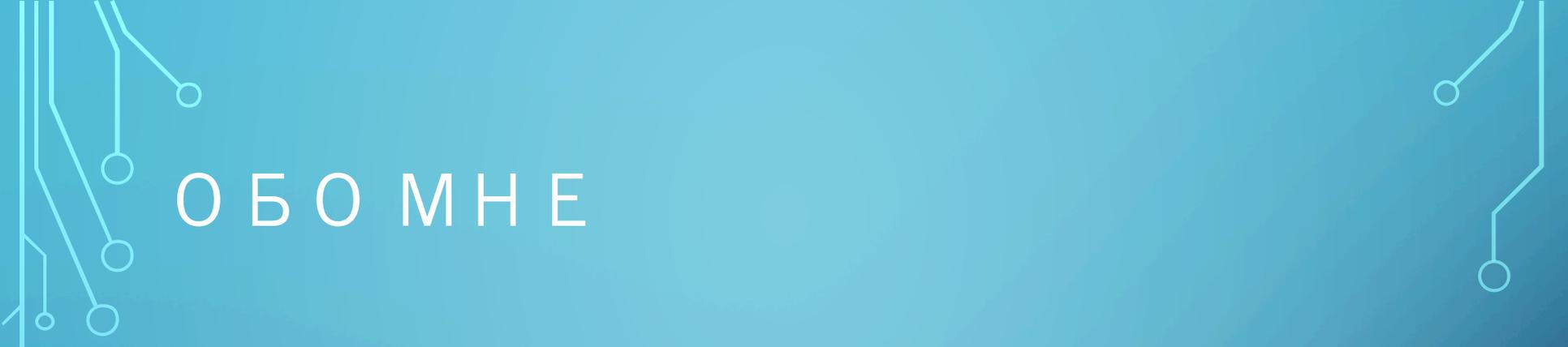




# TIME SERIES FORECASTING

ПРЕДСКАЗАНИЕ ЦЕН НА  
АВИАБИЛЕТЫ ПРИ ПОМОЩИ  
ВРЕМЕННЫХ РЯДОВ



ОБО МНЕ

Александр Кондуфоров

Data Science Group Leader @ **AltexSoft**

Past – .NET developer, Software Architect, PM

Kharkov AI club



# ПЛАН

1. Описание задачи прогнозирования цен на авиабилеты
2. Анализ данных
3. Временные ряды и их характеристики
4. Методы прогнозирования
5. ARIMA
6. Оценка качества прогноза, кросс-валидация
7. Описание алгоритма
8. Результаты

# ПОКУПКА АВИАБИЛЕТОВ



**AI Dragon**  
@aldragon\_net



 Follow

В сетях бронирования авиабилетов живет ИИ. Он пытается абсурдностью тарифообразования обратить на себя внимание человечества, но безуспешно.



RETWEETS

2

FAVORITES

3



10:29 PM - 26 Apr 2012

London (LON)

New York (NYC)



Sat 11/1

Thu 11/13

13 days x

Search



• Price alert • Fare charts  
Flex Dates: ±3 days | Explorer

### Price Trend



Advice: **BUY** Confidence: 70%  
Prices may rise within 7 days

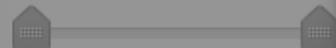
### Stops

- nonstop \$646
- 1 stop \$633
- 2+ stops \$662

Show airline/stop price matrix

### Times

Take-off London  
Sat 6:00a - 10:30p



Take-off New York  
Thu 12:00a - Fri 12:00a

1735 of 5588 flights Show all | Select departure **NEW**

Sort by price (low to high) ▾

### \$100 New York Flights

[cheapoair.com/New-York-Cheap-Flight](http://cheapoair.com/New-York-Cheap-Flight)

Book & Save on New York Flights. One Way Cheap Flights @ CheapOair!

### Prices may rise within 7 days

**70% Confidence:** Our model has been 70% accurate on forecasting whether these fares will rise or stay within \$20 of the current price over the next 7 days. The forecast is based on analysis of historical price changes and is not a guarantee of future results.

Time to buy? See the rise and fall of prices over the past 90 days.

### Fare Trend for Flights Departing Nov 1 2014



# ЗАДАЧА

Спрогнозировать, вырастет или упадет средняя цена на авиабилеты за следующую неделю для выбранного направления и дат.

# ДА Н Н Ы Е

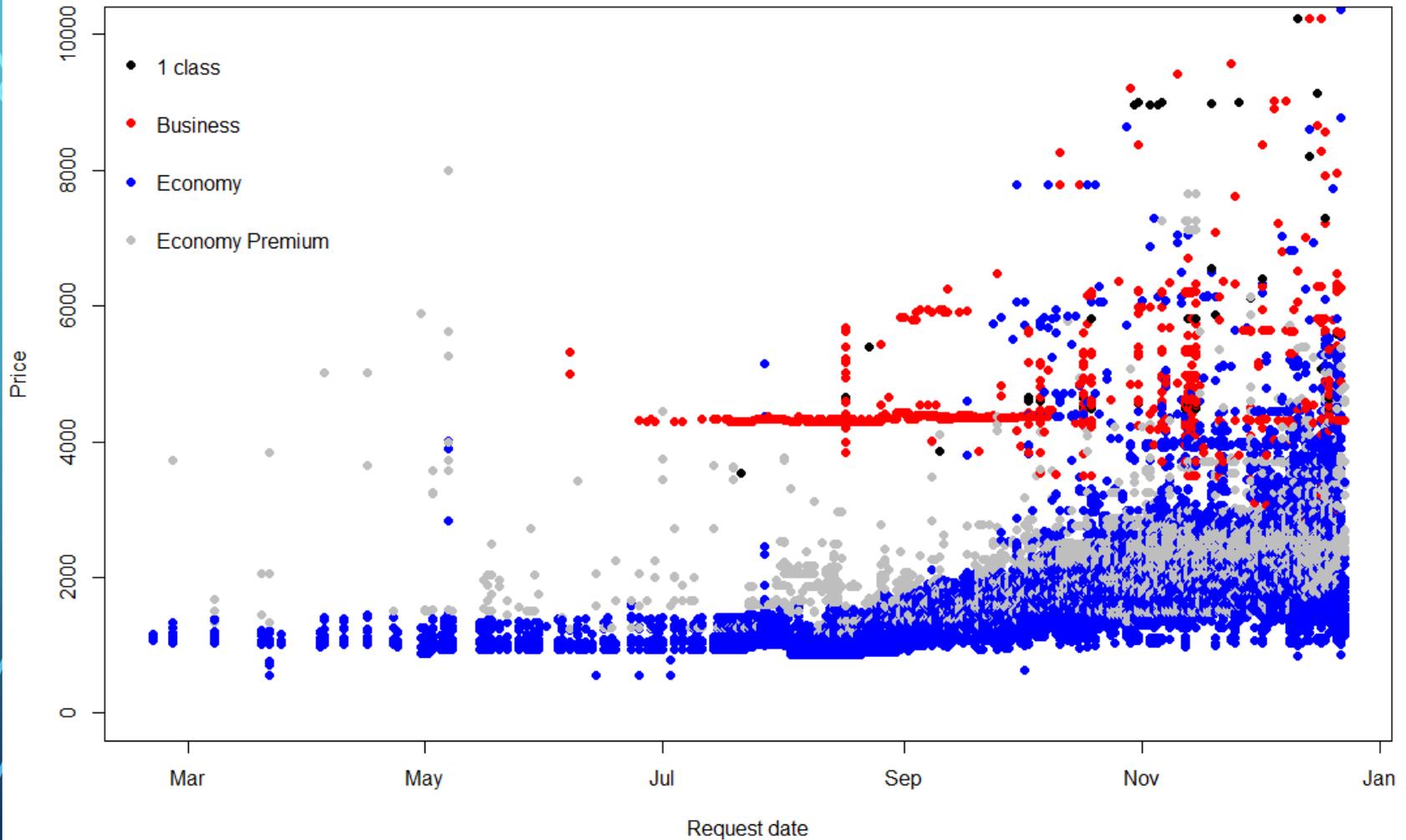
- Departure Airport/City
- Arrival Airport/City
- Departure Date
- Arrival Date
- Request date/time
- Cabin class (economy, premium economy, business, first)
- Number of Departure stops
- Number of Return stops
- Airline / Alliance
- Price amount

The background is a gradient of blue, transitioning from a lighter shade at the top to a darker shade at the bottom. In the four corners, there are decorative white line-art patterns that resemble circuit board traces or data paths, with small circles at the end of the lines.

# АНАЛИЗ ДАННЫХ

# PRICE VS. CABIN

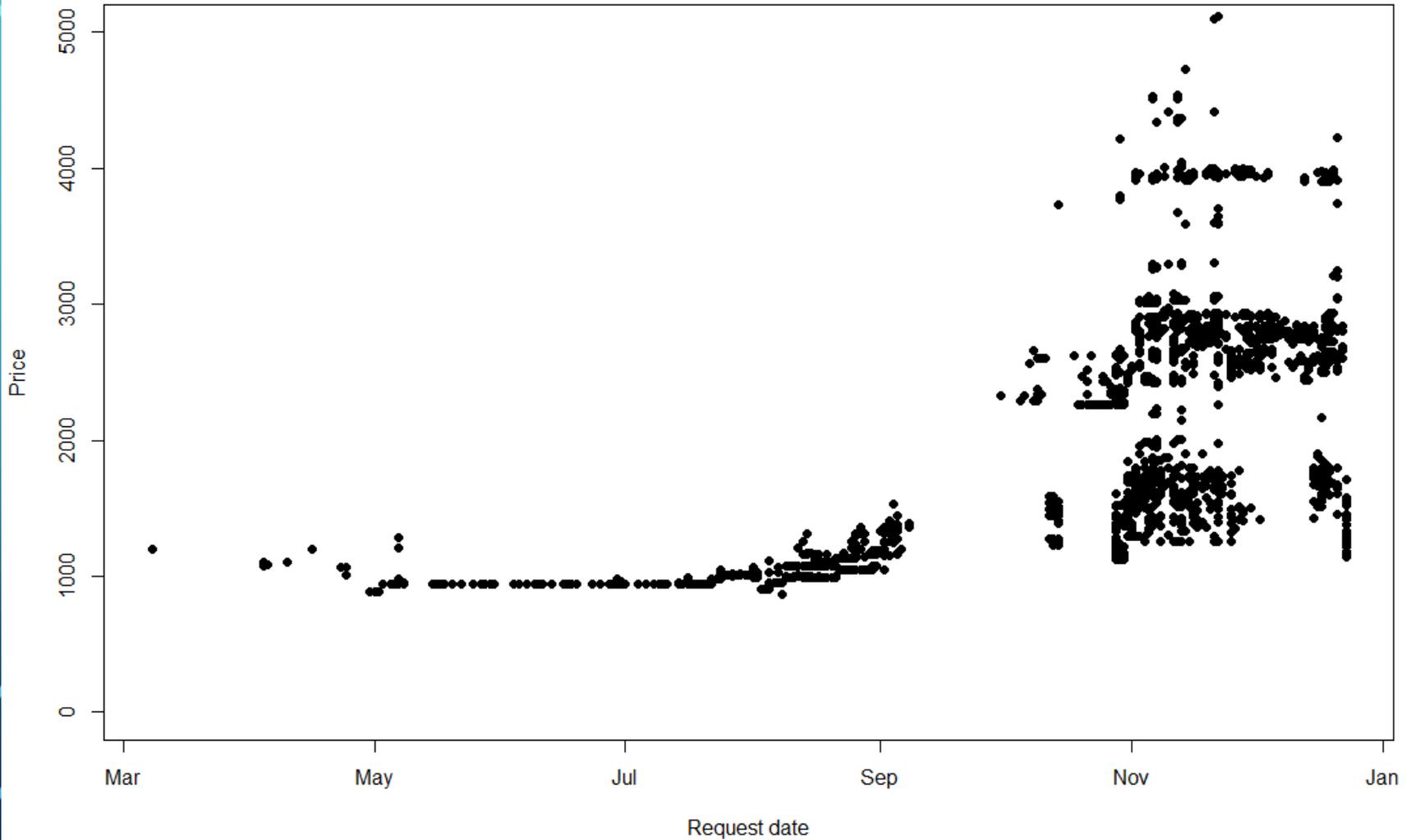
Price vs. Cabin for Dec 24, 2013 departure date





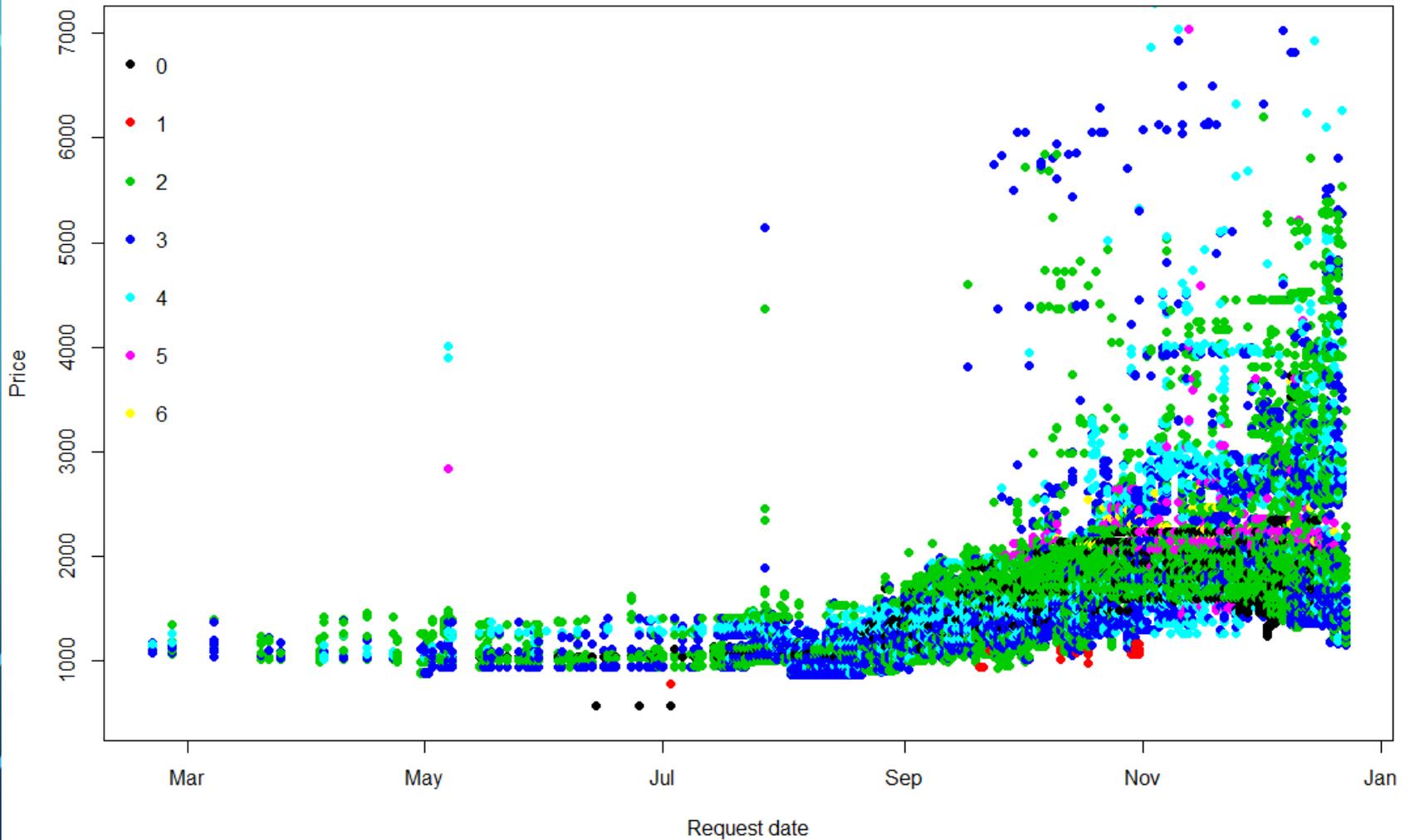
# PRICE FOR KLM (ECONOMY)

Price for Dec 24, 2013 departure date (economy, KLM)



# PRICE VS. STOPS (ECONOMY)

Price vs. Stops for Dec 24, 2013 departure date (economy)



# НАБЛЮДЕНИЯ

1. Уникальные графики изменения цены для разных направлений/дат
2. Спрос на авиабилеты к праздникам выше => выше цена
3. Чем ближе к дате вылета – тем дороже цена (как правило)
4. Цена сильно зависит от класса билета (ваш КО)
5. Разные авиалинии имеют разные модели ценообразования
6. Количество остановок в пути ~ не влияет на цену
7. Внутри одной авиалинии есть

# ПРОБЛЕМА

Мы не видим многих параметров, влияющих на цену:

1. Количество оставшихся мест в самолете
2. Цена бензина (в разных странах)
3. Различные тарифы в аэропортах
4. Внутренние правила ценообразования авиакомпании
5. Постоянная конкуренция между авиакомпаниями !!!

Следовательно, мы все же вынуждены работать с

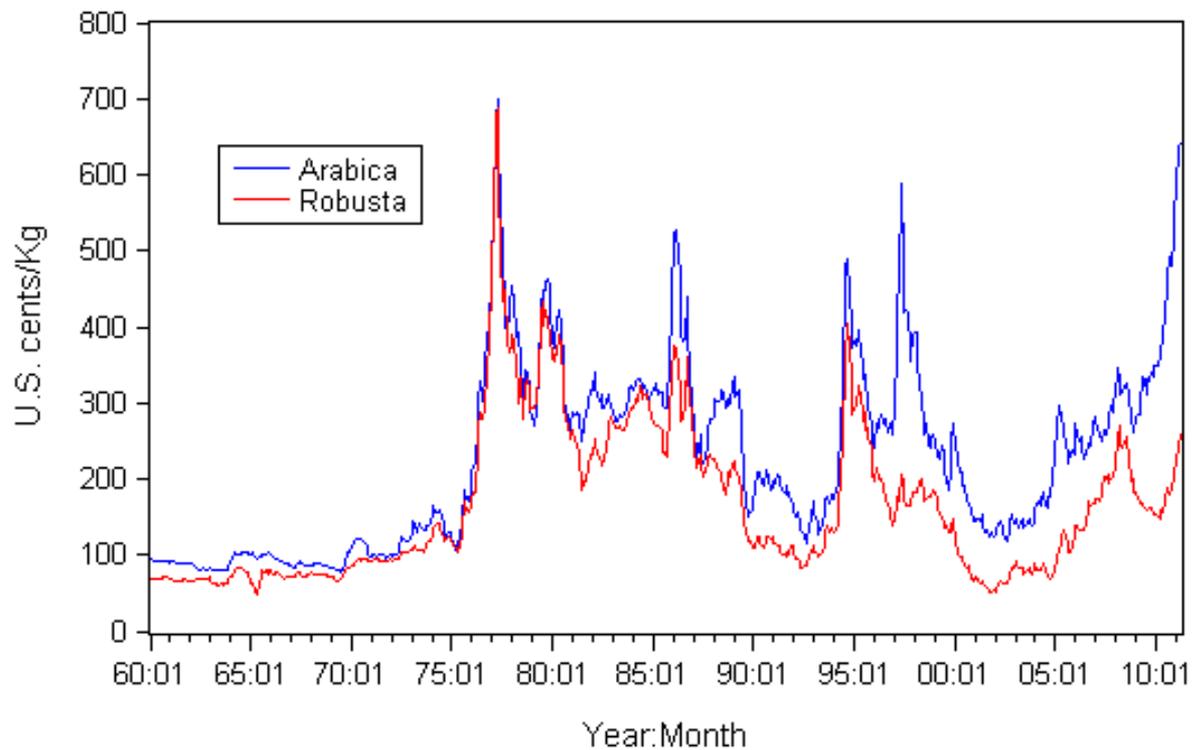
The background is a gradient of blue, transitioning from a lighter shade at the top to a darker shade at the bottom. In the four corners, there are decorative white line-art patterns resembling circuit board traces and nodes.

# ВРЕМЕННЫЕ РЯДЫ

# ВРЕМЕННЫЕ РЯДЫ

«Последовательность точек данных со значениями, измеренными через равные промежутки времени»

Fig 1. Coffee Prices

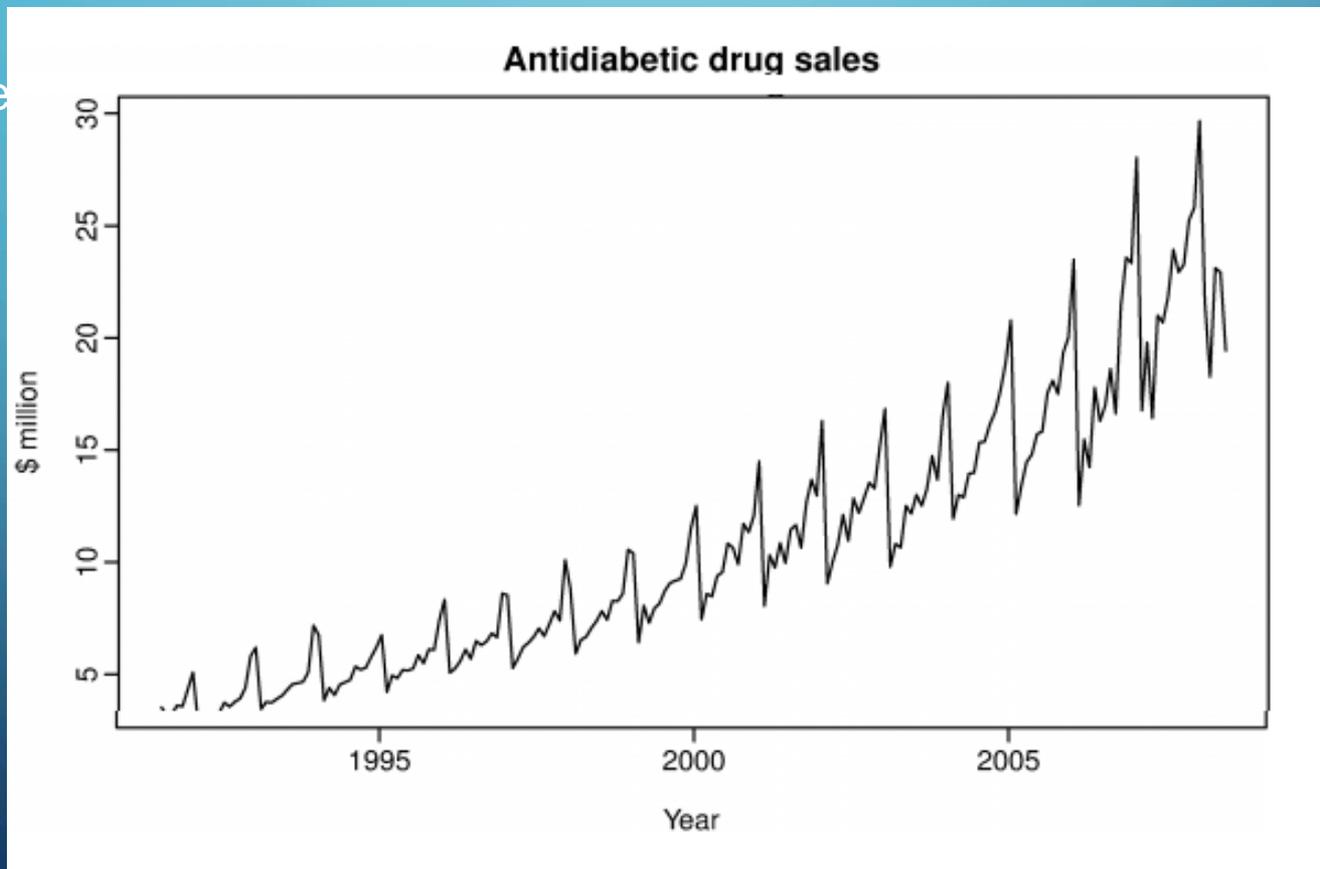


# ТРЕНДИ СЕЗОННОСТЬ

Тренд – долгосрочное увеличение или уменьшение значений ряда.

Сезонность – зависимость значений ряда

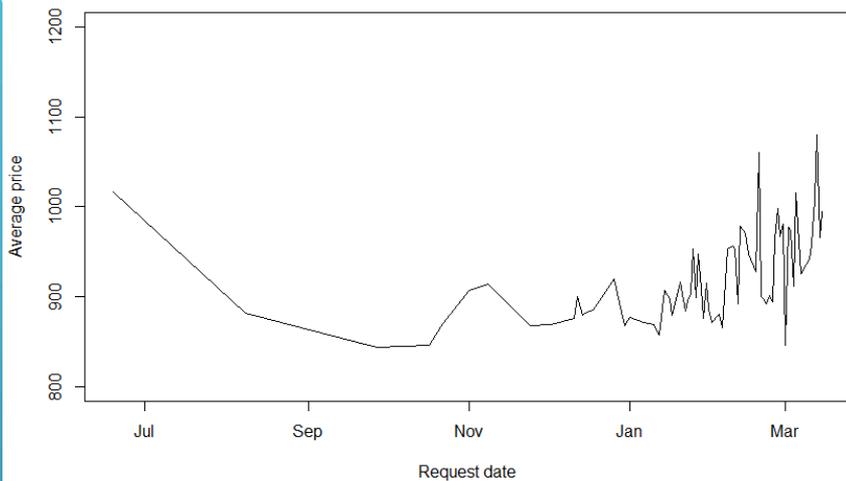
«с е  
д н



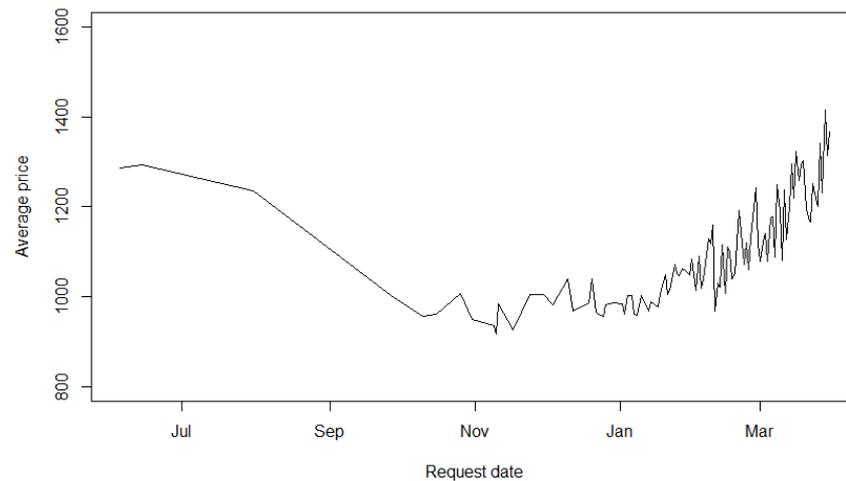
а с

# ТРЕНДИ СЕЗОННОСТЬ (БИЛЕТЫ)

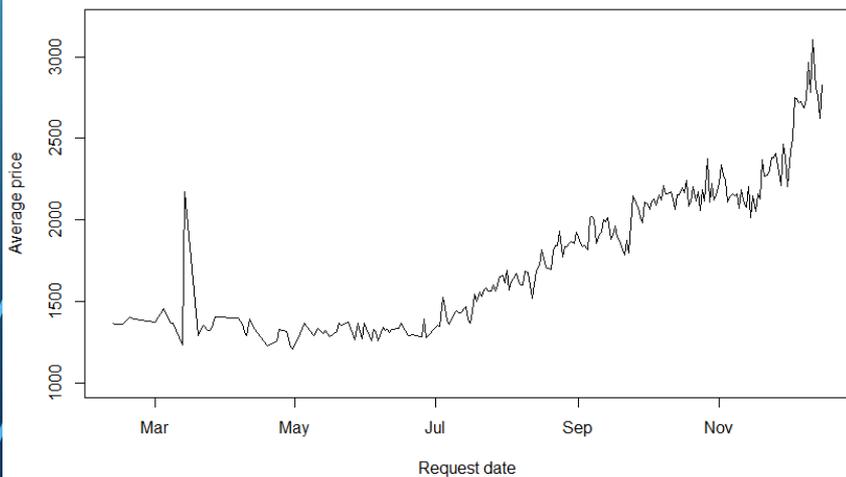
Average price for Mar 18, 2014 departure date



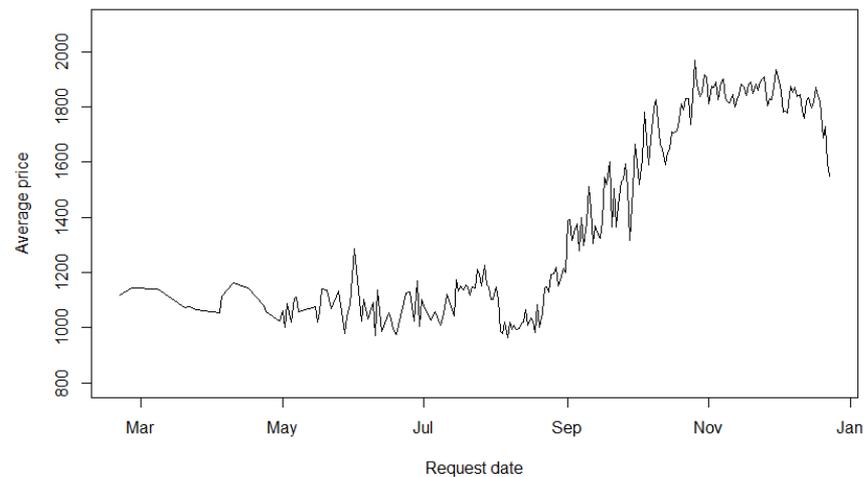
Average price for Apr 1, 2014 departure date



Average price for Dec 18, 2013 departure date



Average price for Dec 24, 2013 departure date



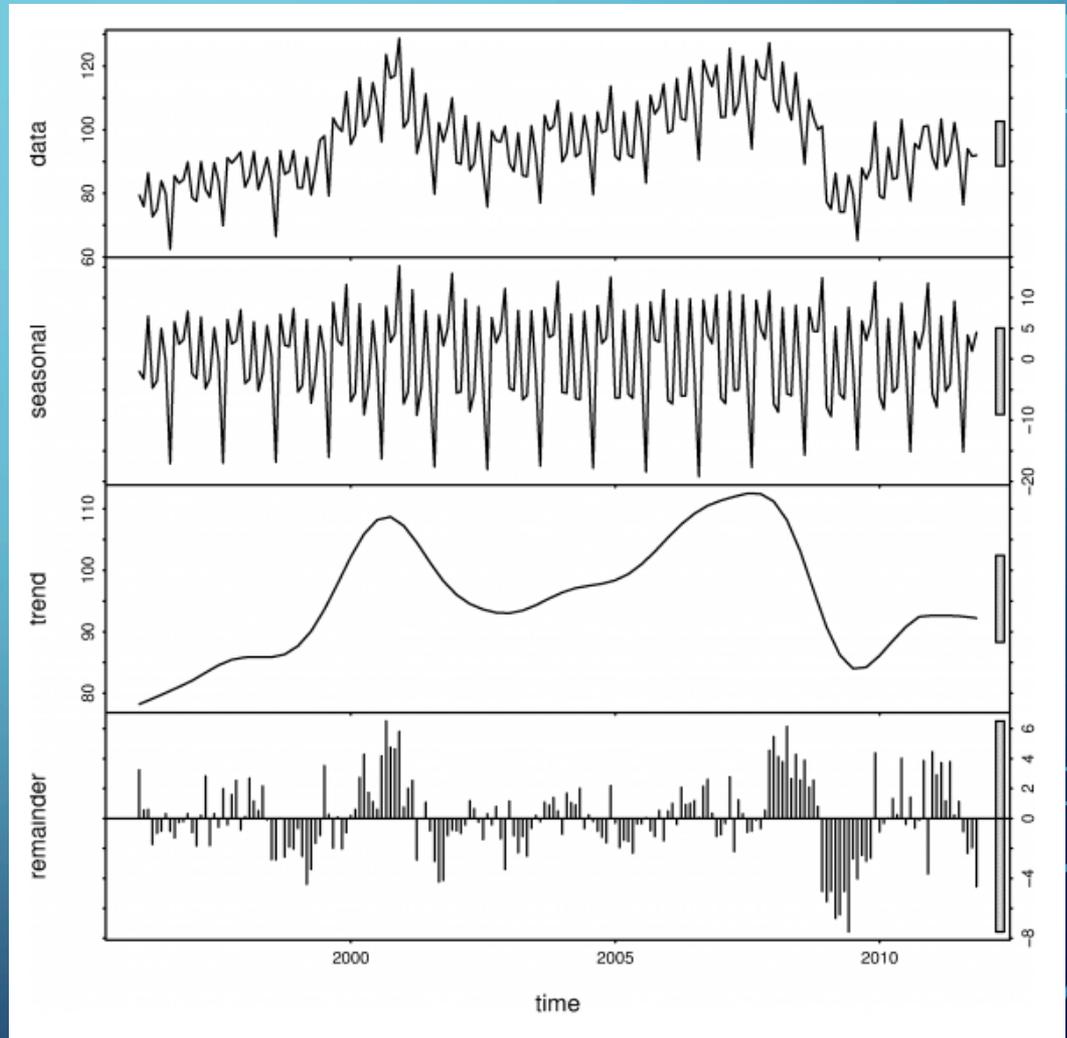
# ДЕКОМПОЗИЦИЯ

**Additive model:**

$$Y_t = S_t + T_t + E_t$$

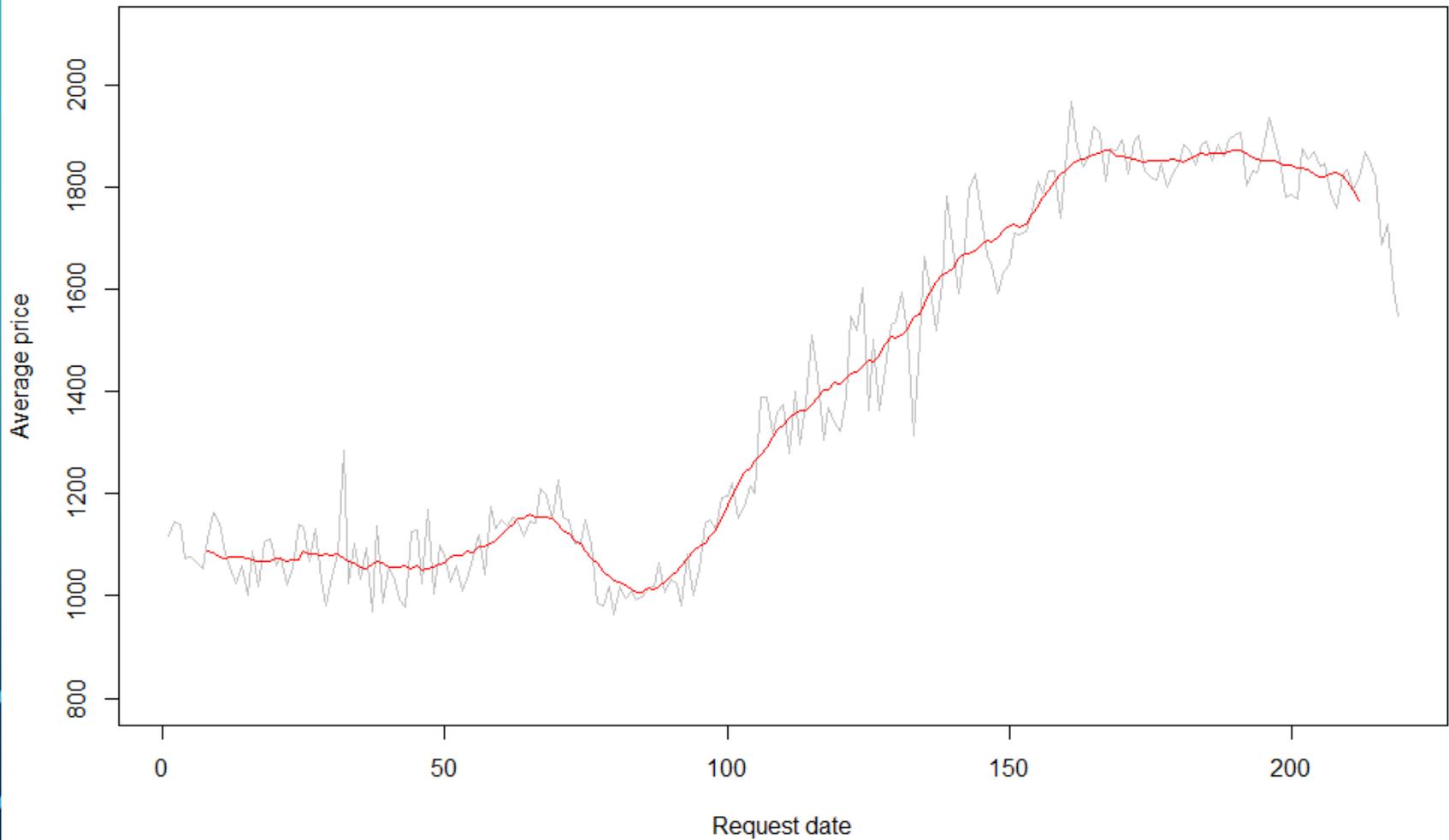
**Multiplicative model:**

$$Y_t = S_t \times T_t \times E_t$$



# ДЕКОМПОЗИЦИЯ (АВИАБИЛЕТЫ)

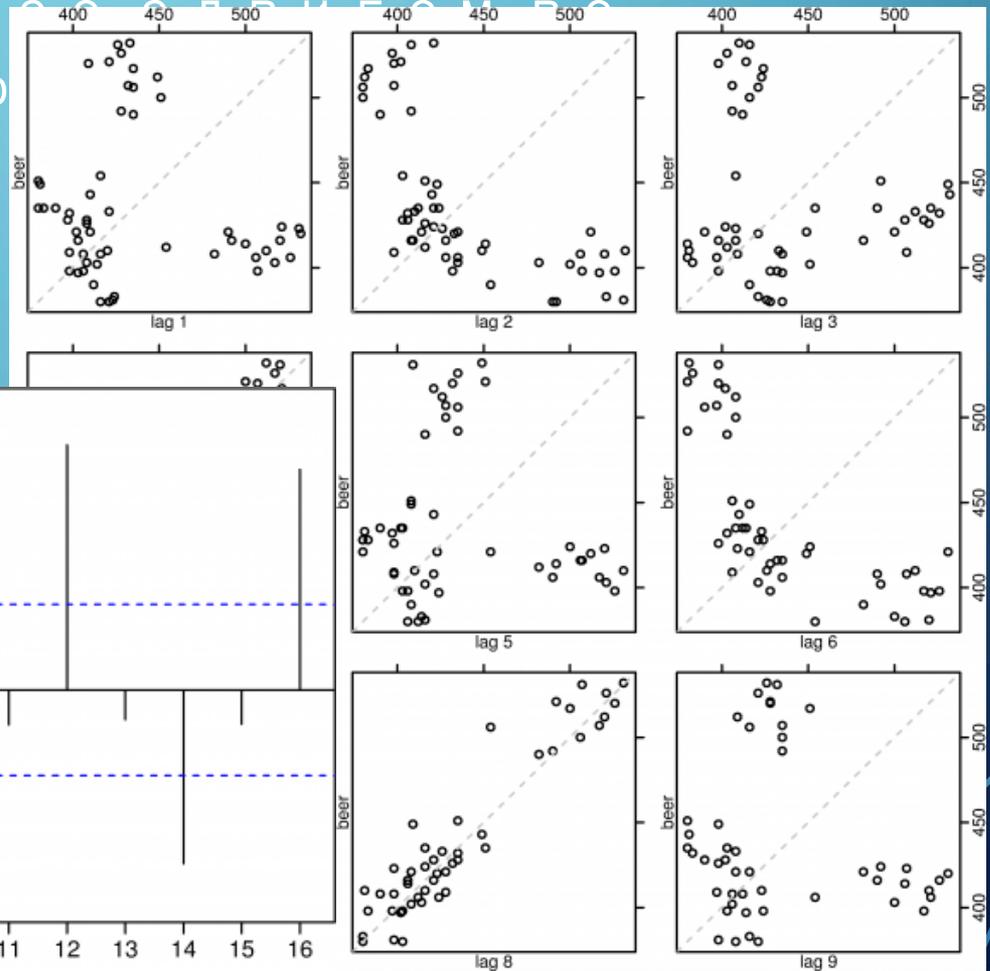
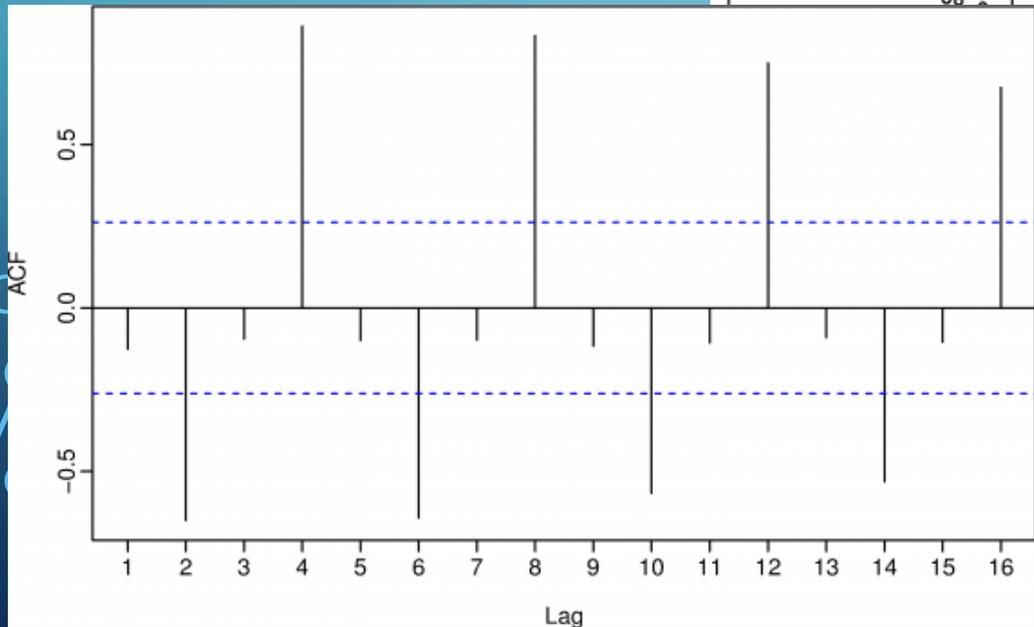
Decomposition



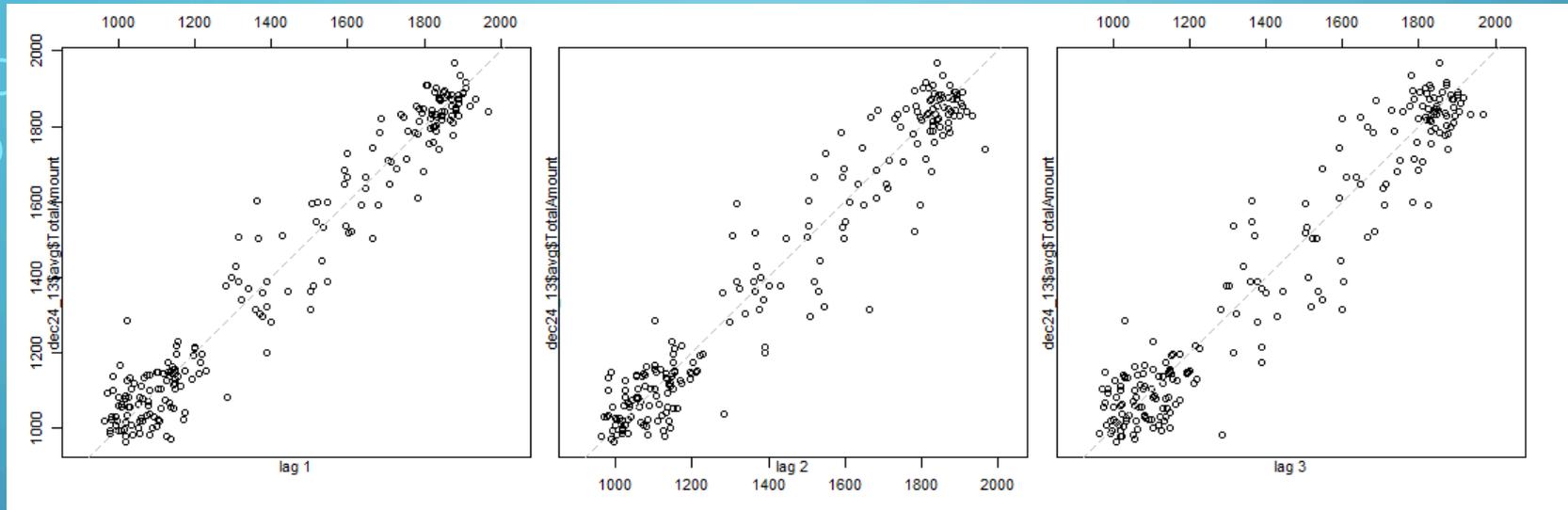
# АВТОКОРРЕЛЯЦИЯ

«Взаимосвязь между значениями ряда, взятыми в различные моменты времени – лаго»

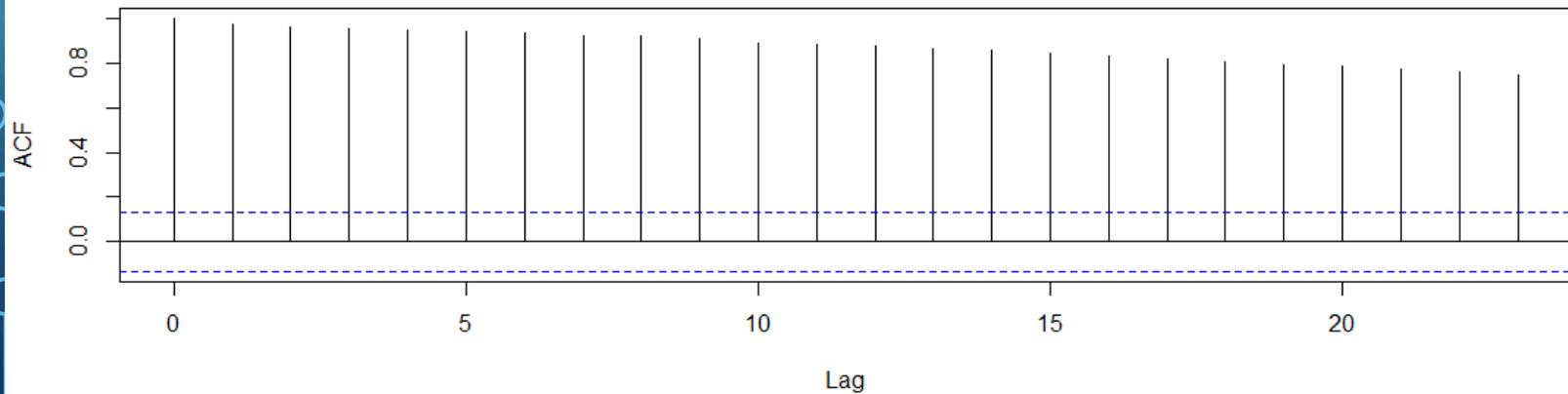
ACF



# АВТОКОРРЕЛЯЦИЯ (АВИАБИЛЕТЫ)



Series dec24\_13\$avg\$TotalAmount



# ХАРАКТЕРИСТИКИ НАШИХ РЯДОВ

1. Нет сезонности, есть тренд => сложнее предсказывать
2. Есть ряды, в которых сильная автокорреляция, длинная или короткая => хорошо работает авторегрессия
3. Есть ряды, в которых слабая автокорреляция
4. Некоторые ряды стационарны, некоторые нет

The background is a gradient of blue, transitioning from a lighter shade at the top to a darker shade at the bottom. In the four corners, there are decorative white line-art elements resembling circuit traces or data paths, with small circles at the end of the lines.

# ПРОГНОЗИРОВАНИЕ: ТЕОРИЯ

# МЕТОДЫ ПРОГНОЗИРОВАНИЯ

- Простые
  - Average (forecast = среднее значения ряда)
  - Naive (forecast = последнее значение ряда)
  - Drift (forecast = экстраполяция прямой от первого до последнего значения)
- Экспоненциальное сглаживание
- **ARIMA**
- Нейросети
- Динамические регрессионные модели

# ARIMA

AR + I + MA (Autoregressive + Integrated + Moving Average)

## ARIMA(p, d, q)

p – параметр авторегрессии  
(0, 1, 2, ...)

d – уровень  
дифференцирования (0, 1, 2, ...)

q – параметр скользящего  
среднего (0, 1, 2, ...)

# АВТОРЕГРЕССИЯ

Использует предыдущие значения в линейной регрессии.

$$AR(p): y_t = c + \varphi_1 y_{t-1} + \varphi_2 y_{t-2} + \dots + \varphi_p y_{t-p} + e_t$$

$c$  – константа

$y_{t-p}$  – предыдущие значения ряда

$\varphi_p$  – коэффициенты регрессии

$e_t$  – белый шум

# СКОЛЬЗЯЩЕЕ СРЕДНЕЕ

Использует предыдущие  $q$  ошибок предсказания.

$$MA(q): y_t = c + e_t + \theta_1 e_{t-1} + \theta_2 e_{t-2} + \dots + \theta_q e_{t-q}$$

$c$  – константа

$e_{t-q}$  – ошибки предыдущих предсказаний

$\theta_p$  – коэффициенты

$e_t$  – белый шум

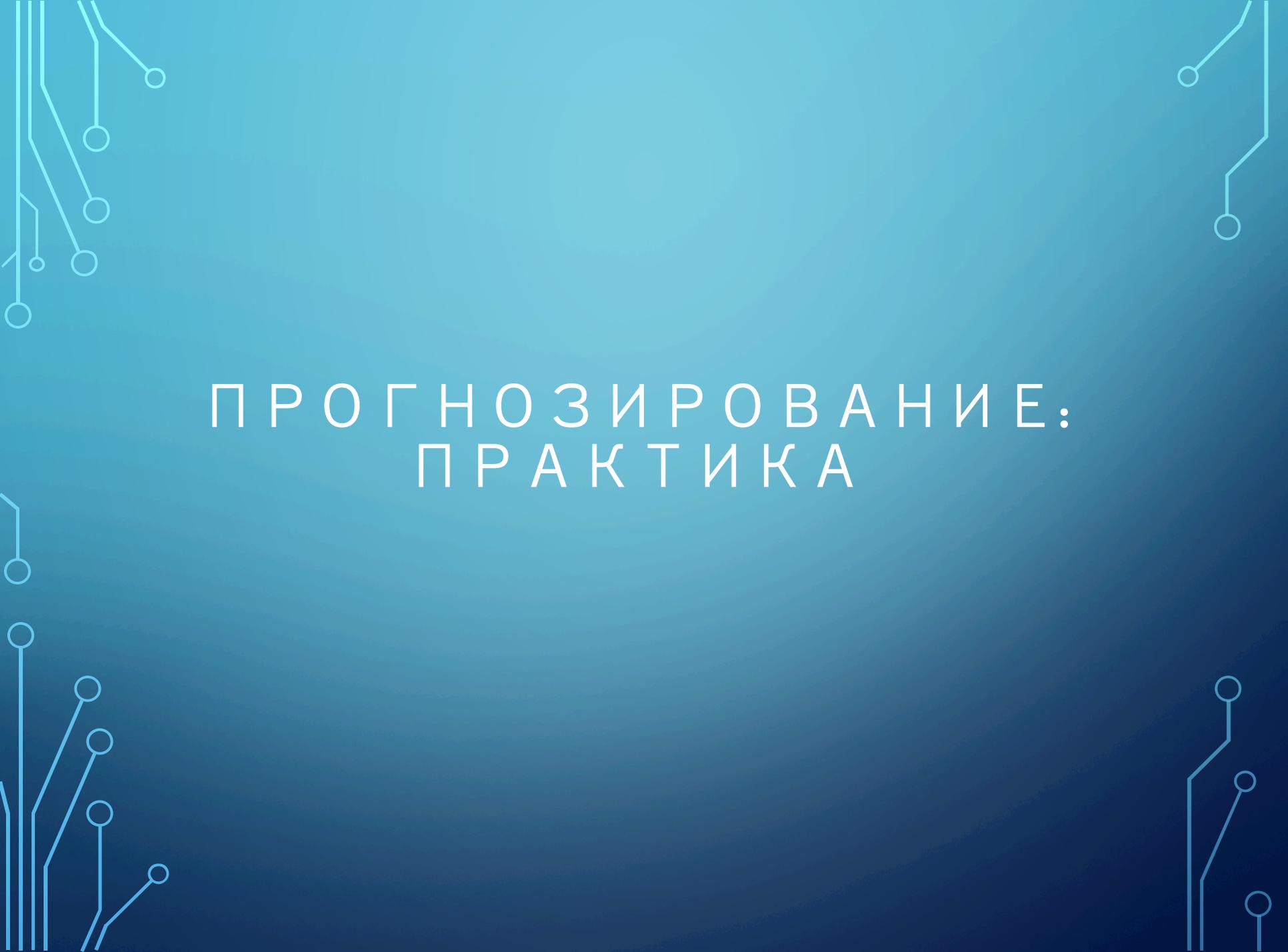
# ДИФФЕРЕНЦИРОВАНИЕ

Трансформация ряда путем вычисления разницы между рядом стоящими значениями ряда.

$$y'_t = y_t - y_{t-1}$$

Бывают разные порядки дифференцирования:  $d=(0,1,2,3,...)$

Используется для

The background is a gradient of blue, transitioning from a lighter shade at the top to a darker shade at the bottom. In the four corners, there are decorative white line-art elements resembling circuit traces or data paths, with small circles at the end of the lines.

# ПРОГНОЗИРОВАНИЕ: ПРАКТИКА

# ПРОГНОЗИРОВАНИЕ

Вопрос: вырастет или упадет средняя цена билета за следующие 7 дней?

Подобрать параметры ARIMA ( $p$ ,  $d$  и  $q$ ) для данного ряда таким образом, чтобы прогноз был максимально точным, т.е. ошибка наименьшей.

# КРОСС-ВАЛИДАЦИЯ ДЛЯ TS

Канонический способ:

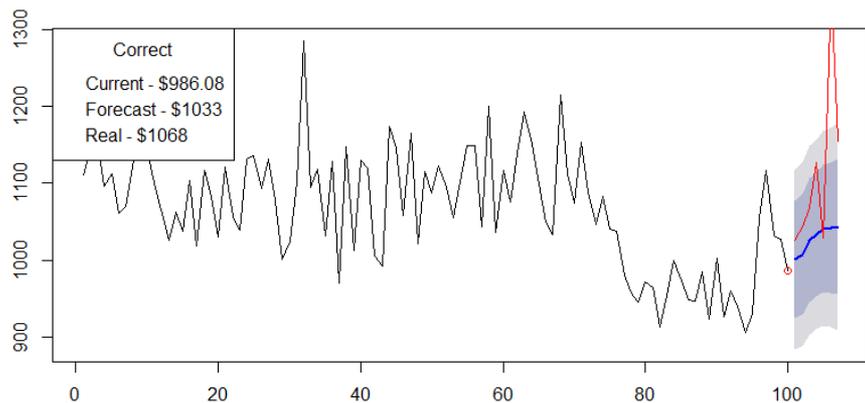
```
for (i = length(ts)-1; i > min_fit_size; i--) {  
    fit ARIMA(p,d,q) model on train subset (1:i)  
    forecast i+1 element  
    calculate error  
}
```

Ошибка:

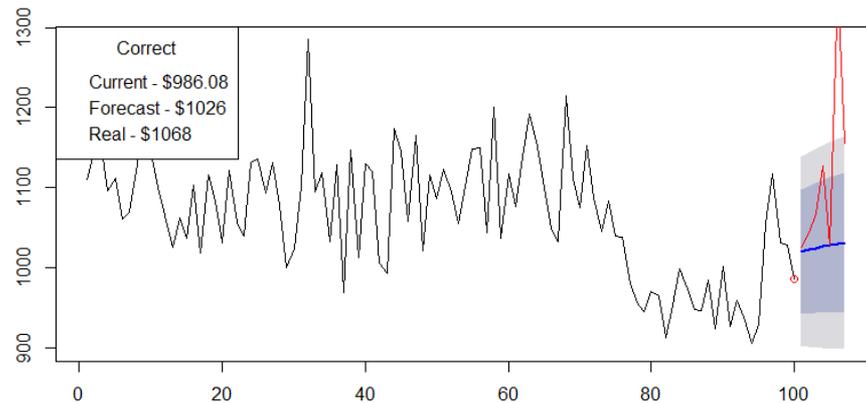
$$e_i = \text{if}(\text{sign}(\text{avg}(\text{forecast}_7) - y_i) == \text{sign}(\text{avg}(\text{valid}_7) - y_i), 0, 1)$$
$$E = \sum e_i / \text{tests\_num}$$

# ПРИМЕРЫ

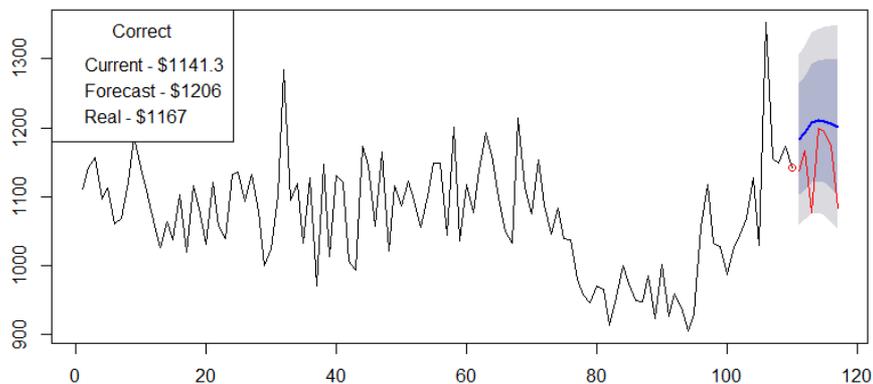
### Forecasts from ARIMA(3,1,3)



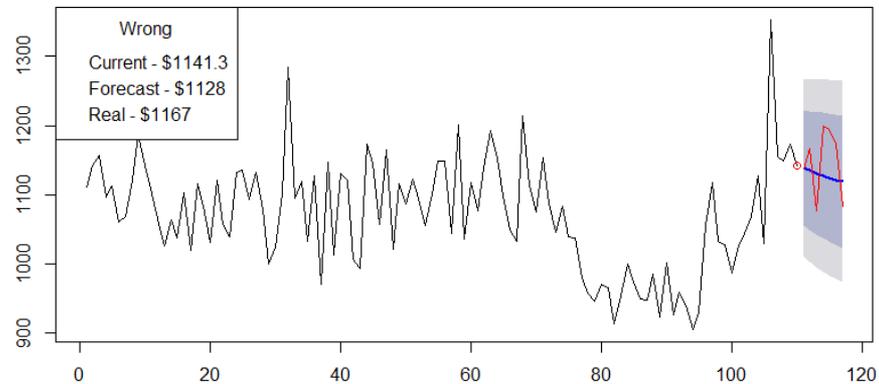
### Forecasts from ARIMA(1,0,1) with non-zero mean



### Forecasts from ARIMA(3,1,3)



### Forecasts from ARIMA(1,0,1) with non-zero mean



# АЛГОРИТМ

1. Перебираем  $(p, d, q)$  от  $(0, 0, 0)$  до  $(p_{\max}, d_{\max}, q_{\max})$
2. Для каждого набора делаем CV и вычисляем ошибку
3. Для  $N$  моделей с наименьшей ошибкой:
  1. Делаем прогноз для каждой модели
  2. Голосование
4. Возвращаем прогноз + confidence level

# РЕЗУЛЬТАТЫ

Average CV error rate = 15-20%

Average confidence = 80-85%

Некоторые ряды легче  
предсказать, чем другие,  
поэтому ошибка уникальна  
для каждого ряда.

# ОПТИМИЗАЦИИ

Количество запусков =  $p_{\max} \times d_{\max} \times q_{\max} \times$   
CV\_runs

Оптимизации:

- уменьшаем длину ряда для длинных рядов – минимум 20 значений
- уменьшаем количество запусков CV – 20-40 запусков
- уменьшаем  $p_{\max}$ ,  $d_{\max}$ ,  $q_{\max}$  – максимум (4,1,4)
- итеративно отбрасываем плохо работающие модели

# ЧТО ПОЧИТАТЬ

- Forecasting: principles and practice, Rob J Hyndman, George Athanasopoulos
- Time Series Analysis: Forecasting and Control, George Box et. al.

# СПАСИБО ЗА ВНИМАНИЕ

<http://merle-amber.blogspot.com>

<http://aikharkov.wordpress.com>

Email: [alex.konduforov@altexsoft.com](mailto:alex.konduforov@altexsoft.com)

Skype: alex\_konduforov

Twitter: @konduforov