

Московский государственный технический университет им. Н.Э. Баумана
Кафедра «Системы обработки информации и управления»



Лабораторная работа №1
по дисциплине
«Методы машинного обучения»
на тему

«Создание "истории о данных" (Data Storytelling)»

Выполнил:
студент группы ИУ5-21М
Чжан Чжибо

Москва — 2021 г.

1. Цель лабораторной работы

Изучение различных методов визуализация данных и создание истории на основе данных.

2. Задание

2.1 Выбрать набор данных (датасет).

2.2 Создать "историю о данных" в виде юпитер-ноутбука, с учетом следующих требований:

История должна содержать не менее 5 шагов (где 5 - рекомендуемое количество шагов). Каждый шаг содержит график и его текстовую интерпретацию.

На каждом шаге наряду с удачным итоговым графиком рекомендуется в юпитер-ноутбуке оставлять результаты предварительных "неудачных" графиков.

Не рекомендуется повторять виды графиков, желательно создать 5 графиков различных видов.

Выбор графиков должен быть обоснован использованием методологии data-to-viz. Рекомендуется учитывать типичные ошибки построения выбранного вида графика по методологии data-to-viz. Если методология Вами отвергается, то просьба обосновать Ваше решение по выбору графика.

История должна содержать итоговые выводы. В реальных "историях о данных" именно эти выводы представляют собой основную ценность для предприятия.

2.3 Сформировать отчет и разместить его в своем репозитории на github.

3. Ход выполнения работы

3.1. Текстовое описание набора данных

Глобальное потепление-это феномен повышения температуры Земли, воздуха и океана за последние пару столетий. Для этого состояния характерна человеческая деятельность, в первую очередь сжигание ископаемого топлива.

В работе использованы наборы данных GlobalTemperatures.csv, GlobalLandTemperatureByCountry.csv, GlobalLandTemperatureByStates.csv из Kaggle, предоставленного пакетом данных Berkeley Earth.

Цель работы: проверка изменения температуры по:

- 1) годам
- 2) месяцам
- 3) временам года
- 4) странам

Самое главное это проверить, правда ли, что температура становится выше.

3.2. Основные характеристики набора данных

Подключим все необходимые библиотеки:

```
In [1]: import pandas as pd
from pandas import DataFrame
from pandas.plotting import scatter_matrix
import matplotlib.pyplot as plt
from matplotlib import rcParams
import plotly.graph_objects as go
import plotly.express as px
from plotly.colors import n_colors
import numpy as np
import datetime as dt
import plotly.express as px
import seaborn as sns
```

Загрузим непосредственно данные:

```
In [6]: #reading the file
GlobalTemp = pd.read_csv("GlobalTemperatures.csv", parse_dates= ['dt'])
GlobalTempCountry = pd.read_csv("GlobalLandTemperaturesByCountry.csv", parse_dates= ['dt'])
GlobalTempState = pd.read_csv("GlobalLandTemperaturesByState.csv", parse_dates= ['dt'])
```

Показание информации о температуре суши и океана на земле за последние несколько лет:

```
In [7]: GlobalTemp.head(5)
```

```
Out[7]:
```

	dt	LandAverageTemperature	LandAverageTemperatureUncertainty	LandMaxTemperature	LandMaxTemperatureUncertainty	LandMinTemperature	LandMinTemperatureUncertainty
0	1750-01-01	3.034	3.574	NaN	NaN	NaN	NaN
1	1750-02-01	3.083	3.702	NaN	NaN	NaN	NaN
2	1750-03-01	5.626	3.076	NaN	NaN	NaN	NaN
3	1750-04-01	8.490	2.451	NaN	NaN	NaN	NaN
4	1750-05-01	11.573	2.072	NaN	NaN	NaN	NaN

Изменение имени столбца для удобства:

```
In [8]: GlobalTemp.rename(columns = ['dt':'Date'], inplace = True)
```

3.3. Визуальное исследование датасета

Найдём время, когда температура начала расти.

```
Year_Temp = GlobalTemp.groupby(GlobalTemp['Date'].dt.year)['LandAverageTemperature', 'LandMaxTemperature',  
                                                             'LandMinTemperature', 'LandAndOceanAverageTemperature'].mean().reset_index()  
Year_Temp.rename(columns = ['Date':'Year'], inplace = True)  
  
fig = go.Figure()  
  
# Add traces  
fig.add_trace(go.Scatter(x=Year_Temp.Year, y=Year_Temp.LandAverageTemperature,  
                        mode='lines',  
                        name='LandAvgTemp',  
                        marker_color='#A9A9A9'))  
fig.add_trace(go.Scatter(x=Year_Temp.Year, y=Year_Temp.LandMaxTemperature,  
                        mode='lines',  
                        name='LandMaxAvgTemp',  
                        marker_color='#BDB76B'))  
fig.add_trace(go.Scatter(x=Year_Temp.Year, y=Year_Temp.LandMinTemperature,  
                        mode='lines',  
                        name='LandMinAvgTemp',  
                        marker_color='#45CE30'))  
  
fig.add_trace(go.Scatter(x=Year_Temp.Year, y=Year_Temp.LandAndOceanAverageTemperature,  
                        mode='lines',  
                        name='Land&OceanAvgTemp',  
                        marker_color='#FFA07A'))  
  
fig.update_layout(  
    height=800,  
    xaxis_title='Years',  
    yaxis_title='Temperatures in degree',  
    title_text='Average Land, Ocean, Minimum, and Maximum Temperatures over the years'  
)  
fig.add_annotation(  
    x=1950,  
    y=2.7,  
    text='1950°')  
fig.add_annotation(  
    x=1972,  
    y=8.4,  
    text='1972°')  
fig.add_annotation(  
    x=1978,  
    y=14.28,  
    text='1978°')  
fig.add_annotation(  
    x=1969,  
    y=15.31,  
    text='1969°')
```

```

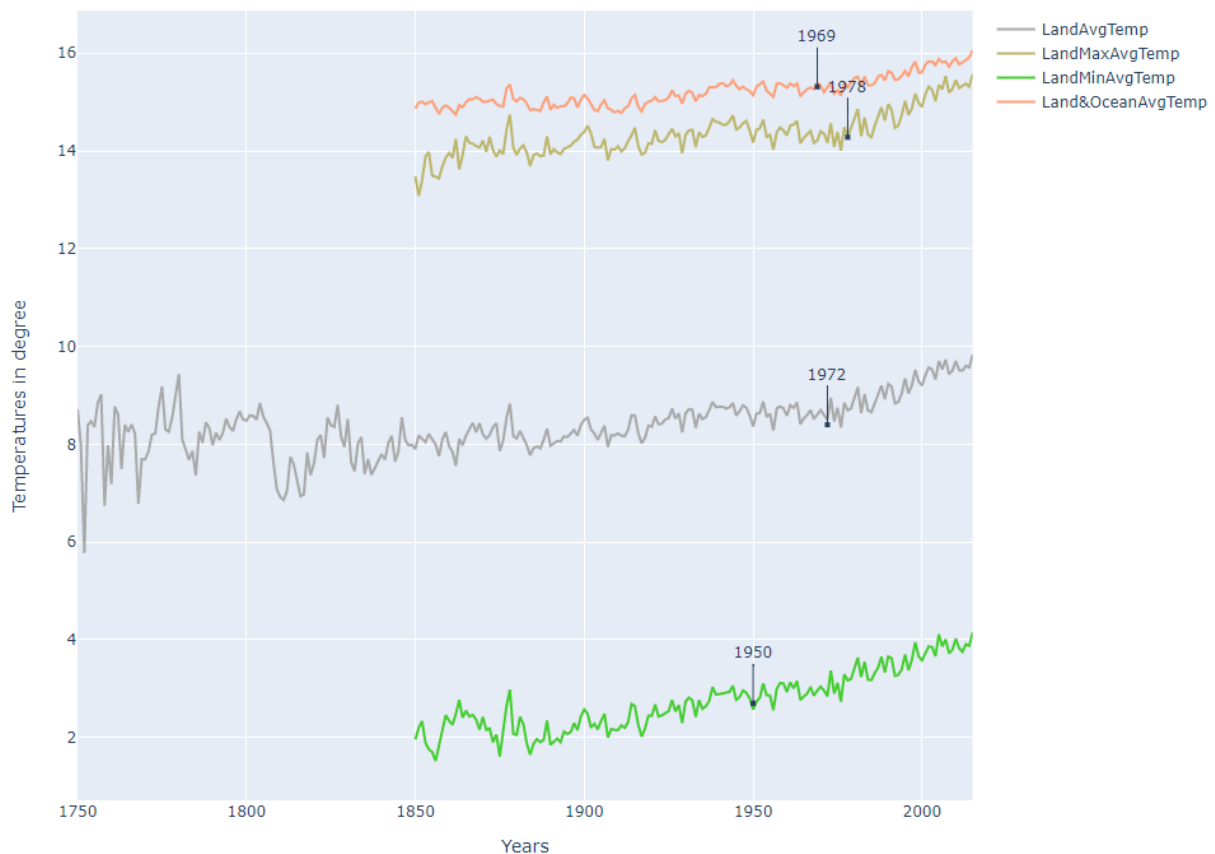
        text="1969")
fig.update_annotations(dict(
    xref="x",
    yref="y",
    showarrow=True,
    arrowhead=7,
    ax=0,
    ay=-40
))

fig.update_layout(showlegend=True)

fig.show()

```

Average Land, Ocean, Minimum, and Maximum Temperatures over the years



Похоже, что данные по LandMaxTemperature, LandMinTemperature, Land&OceanAverageTemperature отсутствуют с 1750 по 1850 год.

Очевидно, что:

Средняя температура земли и океана увеличилась на 1,19 градуса с 1850 по 2015 год.

Максимальная температура земли увеличилась на 2,1 градуса с 1850 по 2015 год.

Средняя температура земли увеличилась на 1,12 градуса с 1750 по 2015 год

Минимальная температура земли достигла максимума в 2,24 градуса с 1850 по 2015 год

Из графика видно, что начиная с 1950 года все параметры стали увеличиваться. Теперь двинемся вперед и проведем ежемесячный анализ, чтобы увидеть, какие месяцы удерживают высокие температуры за эти годы

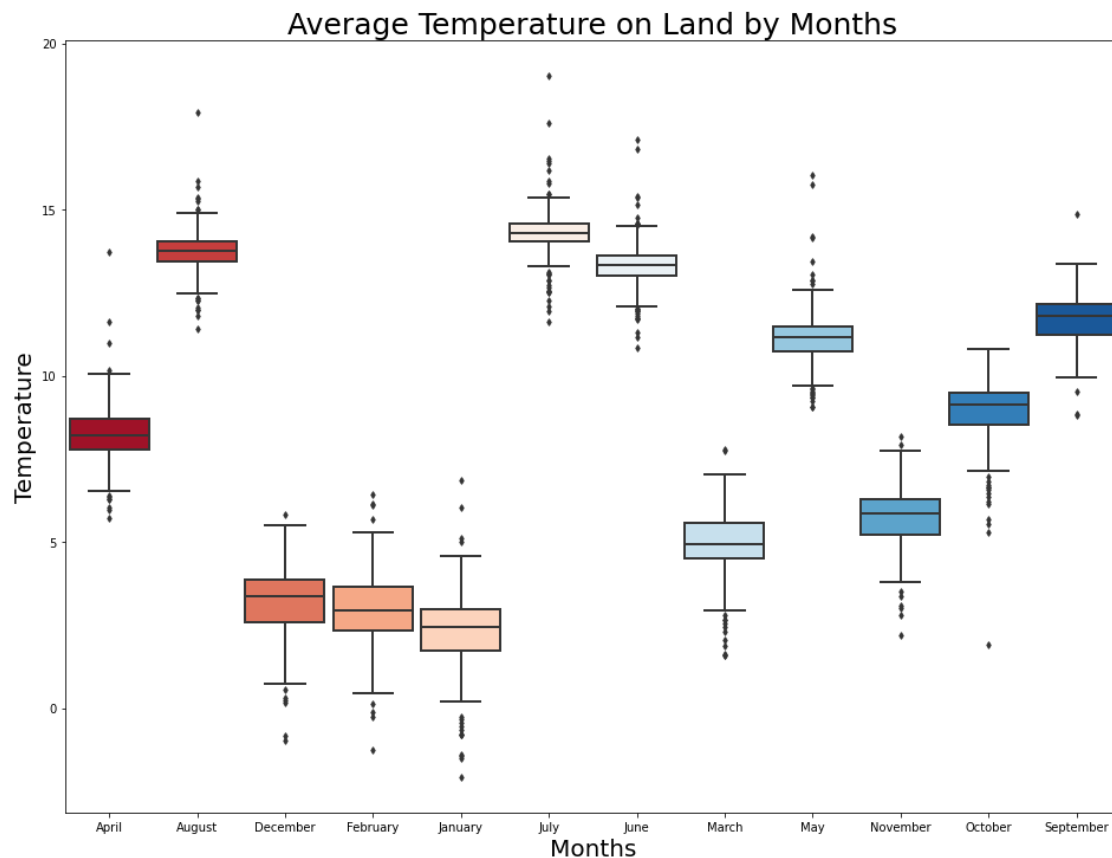
Ежемесячный анализ температуры по годам:

```
In [10]: GlobalTemp["Year"] = pd.DatetimeIndex(GlobalTemp["Date"]).year
GlobalTemp["Month"] = pd.DatetimeIndex(GlobalTemp["Date"]).month
GlobalTemp["Month"] = GlobalTemp["Month"].astype(str)
GlobalTemp.loc[GlobalTemp["Month"]=='1', 'Month'] = 'January'
GlobalTemp.loc[GlobalTemp["Month"]=='2', 'Month'] = 'February'
GlobalTemp.loc[GlobalTemp["Month"]=='3', 'Month'] = 'March'
GlobalTemp.loc[GlobalTemp["Month"]=='4', 'Month'] = 'April'
GlobalTemp.loc[GlobalTemp["Month"]=='5', 'Month'] = 'May'
GlobalTemp.loc[GlobalTemp["Month"]=='6', 'Month'] = 'June'
GlobalTemp.loc[GlobalTemp["Month"]=='7', 'Month'] = 'July'
GlobalTemp.loc[GlobalTemp["Month"]=='8', 'Month'] = 'August'
GlobalTemp.loc[GlobalTemp["Month"]=='9', 'Month'] = 'September'
GlobalTemp.loc[GlobalTemp["Month"]=='10', 'Month'] = 'October'
GlobalTemp.loc[GlobalTemp["Month"]=='11', 'Month'] = 'November'
GlobalTemp.loc[GlobalTemp["Month"]=='12', 'Month'] = 'December'
year_month = GlobalTemp.groupby(by = ['Year', 'Month']).mean().reset_index()
# Figure size
plt.figure(figsize=(16,12))

# The plot
sns.boxplot(x = 'Month', y = 'LandAverageTemperature', data = year_month, palette = "RdBu", saturation = 1, width = 0.9, fliersize=4, linewidth=1)

# Make pretty
plt.title('Average Temperature on Land by Months', fontsize = 25)
plt.xlabel('Months', fontsize = 20)
plt.ylabel('Temperature', fontsize = 20)
```

Out[10]: Text(0, 0.5, 'Temperature')

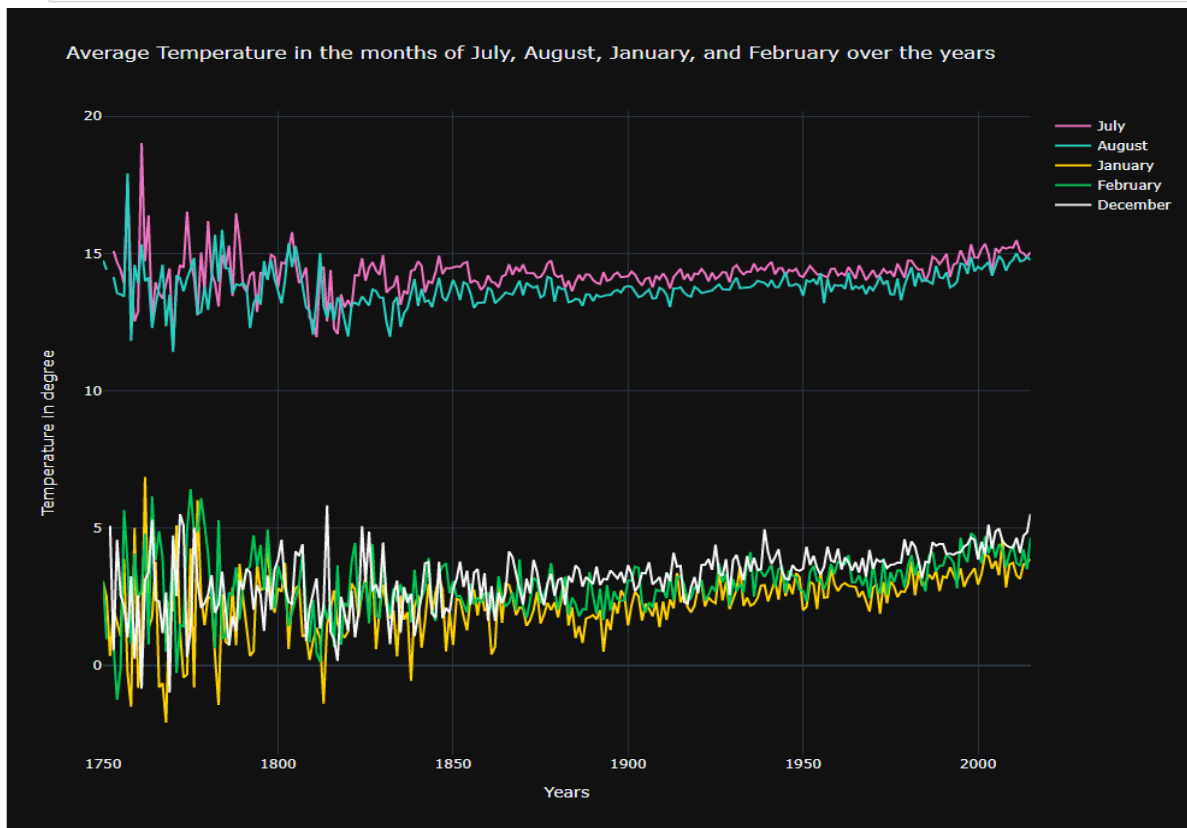


Как мы знаем, температура, по-видимому, была высокой в июле и августе и низкой в январе, феврале и декабре на протяжении 250 лет. Дальше копнем глубже и увидим повышение температуры с годами за эти месяцы

```
In [11]: month_temp = GlobalTemp.groupby(by = ['Year', 'Month']).mean().reset_index()

July = month_temp.loc[month_temp['Month'] == 'July', :]
August = month_temp.loc[month_temp['Month'] == 'August', :]
January = month_temp.loc[month_temp['Month'] == 'January', :]
February = month_temp.loc[month_temp['Month'] == 'February', :]
December = month_temp.loc[month_temp['Month'] == 'December', :]
fig1 = go.Figure()
for template in ['plotly_dark']:
    fig1.add_trace(go.Scatter(x=July['Year'], y=July['LandAverageTemperature'],
                             mode='lines',
                             name='July',
                             marker_color='#f075c2'))
    fig1.add_trace(go.Scatter(x=August['Year'], y=August['LandAverageTemperature'],
                              mode='lines',
                              name='August',
                              marker_color='#28d2c2'))
    fig1.add_trace(go.Scatter(x=January['Year'], y=January['LandAverageTemperature'],
                              mode='lines',
                              name='January',
                              marker_color='#ffd201'))
    fig1.add_trace(go.Scatter(x=February['Year'], y=February['LandAverageTemperature'],
                              mode='lines',
                              name='February',
                              marker_color='#00c957'))
    fig1.add_trace(go.Scatter(x=December['Year'], y=December['LandAverageTemperature'],
                              mode='lines',
                              name='December',
                              marker_color='#ff7f7f'))
    fig1.update_layout(
        height=800,
        xaxis_title='Years',
        yaxis_title='Temperature in degree',
        title_text='Average Temperature in the months of July, August, January, and February over the years',
        template=template)

fig1.show()
```



Видно, что:

Средняя температура поднялась на 0,041 градуса в июле месяце с 1753 по 2015 год.

Средняя температура поднялась на 0,005 градуса в августе месяце с 1750 по 2015 год.

Средняя температура поднялась почти на 1 градус в январе с 1750 по 2015 год.

Средняя температура поднялась на 1,58 градуса в феврале месяце с 1750 по 2015 год.

Средняя температура поднялась на 2,746 градуса в декабре с 1750 по 2015 год.

По разнице температур вышеописанных месяцев можно сказать, что в холодные месяцы разница температур выше, чем в жаркие. Это признак того, что холодные месяцы становятся все жарче.

Теперь проверим, не становится ли какое-нибудь время года жарче с годами. Посмотрим сезонное повышение температуры за данные годы. Поскольку мы имеем дело с данными временных рядов, будет предпочтительнее снова использовать линейный график

```
In [12]: month_season = {
    "January": "Winter",
    "February": "Winter",
    "March": "Spring",
    "April": "Spring",
    "May": "Spring",
    "June": "Summer",
    "July": "Summer",
    "August": "Summer",
    "September": "Autumn",
    "October": "Autumn",
    "November": "Autumn",
    "December": "Winter"
}

GlobalTemp['Season'] = ''

for month, season in month_season.items():
    GlobalTemp.loc[GlobalTemp['Month'] == month, 'Season'] = season

year_season = GlobalTemp.groupby(by=['Year', 'Season']).mean().reset_index()

Winter = year_season.loc[year_season['Season'] == 'Winter', :]
Spring = year_season.loc[year_season['Season'] == 'Spring', :]
Summer = year_season.loc[year_season['Season'] == 'Summer', :]
Autumn = year_season.loc[year_season['Season'] == 'Autumn', :]

fig2 = go.Figure()
for template in ['plotly_white']:
    fig2.add_trace(go.Scatter(x=Winter['Year'], y=Winter['LandAverageTemperature'],
                             mode='lines',
                             name='Winter',
                             marker_color='#338B8B'))
    fig2.add_trace(go.Scatter(x=Spring['Year'], y=Spring['LandAverageTemperature'],
                             mode='lines',
                             name='Spring',
                             marker_color='#FFB5C5'))
    fig2.add_trace(go.Scatter(x=Summer['Year'], y=Summer['LandAverageTemperature'],
                             mode='lines',
                             name='Summer',
                             marker_color='#87CEFF'))
    fig2.add_trace(go.Scatter(x=Autumn['Year'], y=Autumn['LandAverageTemperature'],
                             mode='lines',
                             name='Autumn',
                             marker_color='#FF8000'))
```



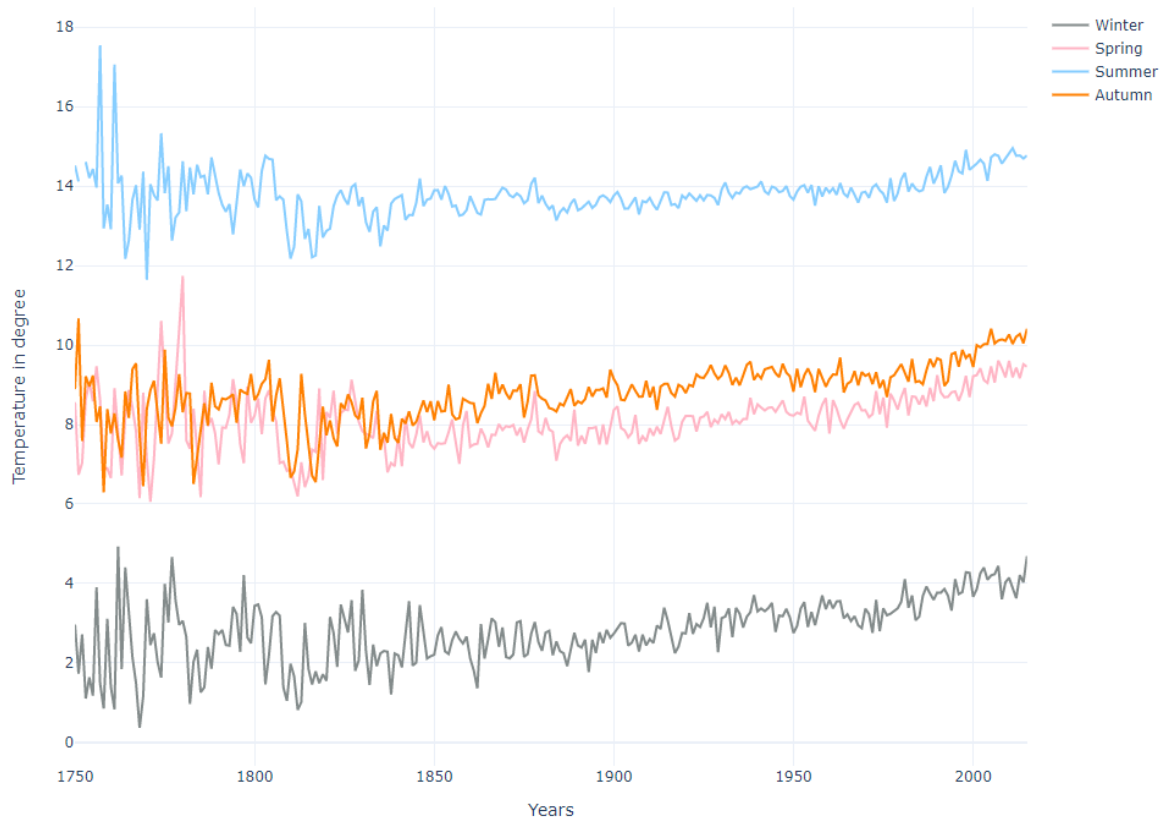
```

        marker_color='#FF8000'))
fig2.update_layout(
    height=800,
    xaxis_title="Years",
    yaxis_title='Temperature in degree',
    title_text='Average Temperature seasonwise over the years',
    template=template)

fig2.show()

```

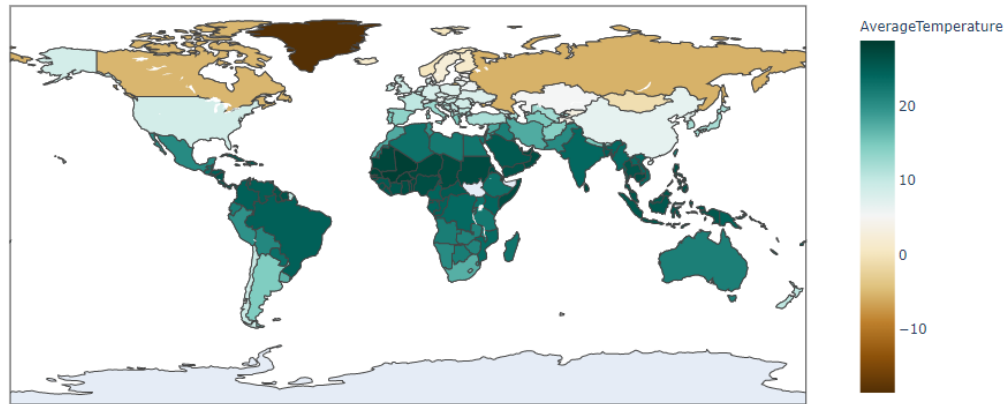
Average Temperature seasonwise over the years



Из графика видно, что весна, лето и зима с каждым годом становятся все жарче. Особенно зимний сезон имеет самый высокий всплеск за последние годы. Таким образом, мы можем сказать, что холодные сезоны становятся все жарче. Перейдем к другому датсету и проверим пострановой средней температуры страны за эти годы:

```
In [13]: country_temp = GlobalTempCountry.groupby(by = ['Country']).mean().reset_index()
fig3 = px.choropleth(country_temp, locations="Country", locationmode = "country names", color="AverageTemperature",
                    color_continuous_scale=px.colors.diverging.BrBG,
                    title="Average Temperature Contrywise Worldwide")
fig3.show()
```

Average Temperature Contrywise Worldwide



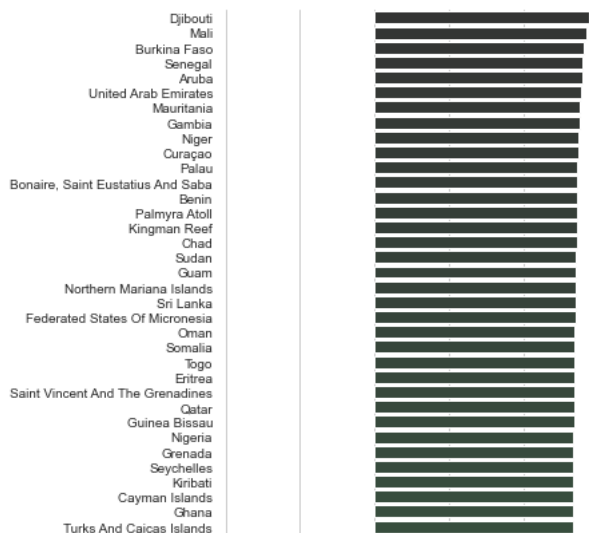
Из приведенного выше хороплета видно, что страны Гренландия, Россия, Канада имеют самую низкую среднюю температуру, а страны африканского континента, такие как Джибути, Мали, Буркина-Фасо, имеют самую высокую среднюю температуру в течение многих лет.

```
In [14]: country_temp_asc = GlobalTempCountry.groupby(by = ['Country']).mean().reset_index().sort_values('AverageTemperature', ascending=False).reset_index()
sns.set(style="whitegrid", font_scale=0.9)

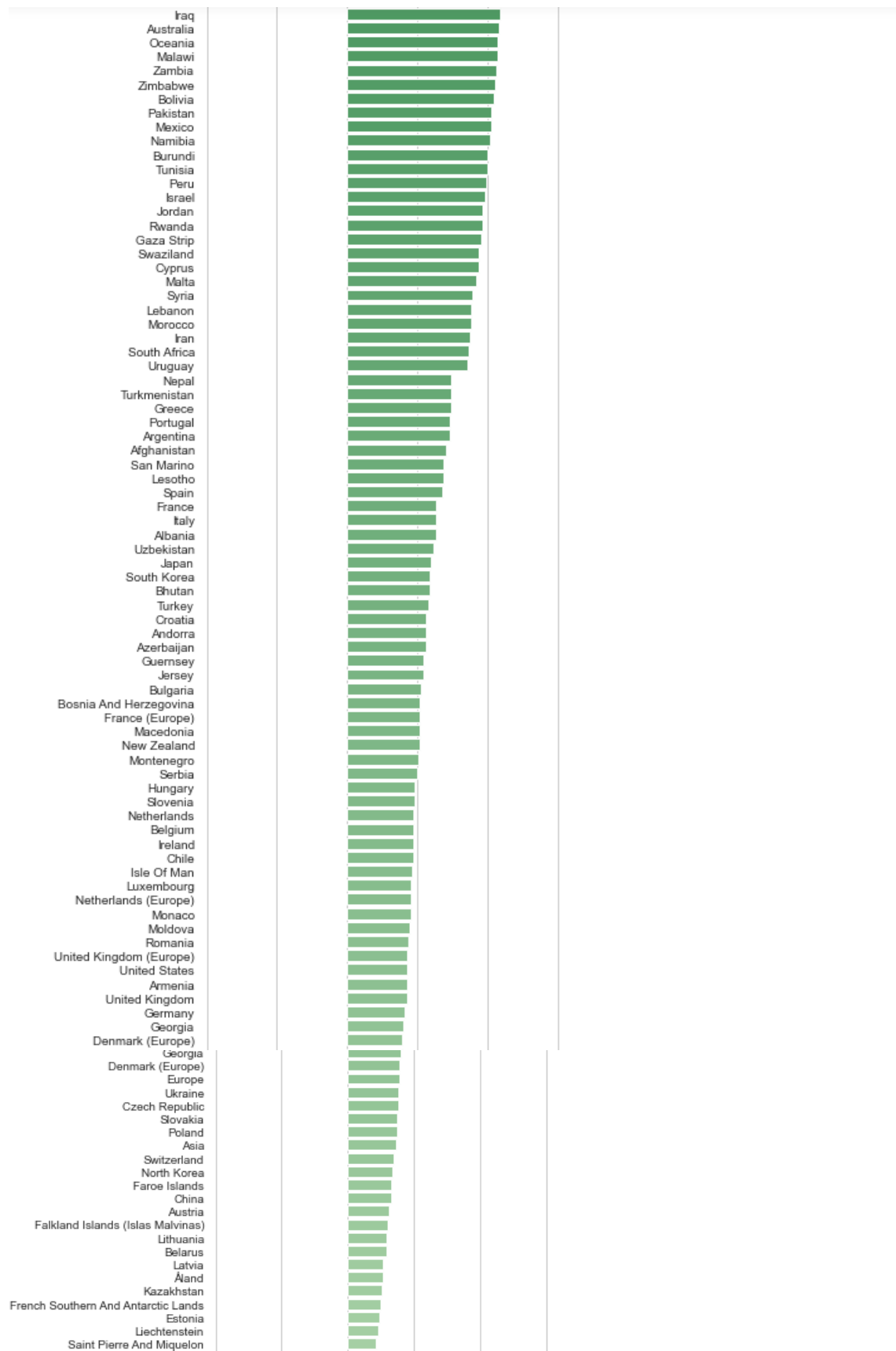
# Initialise the matplotlib figure
f, ax = plt.subplots(figsize=(6, 50))

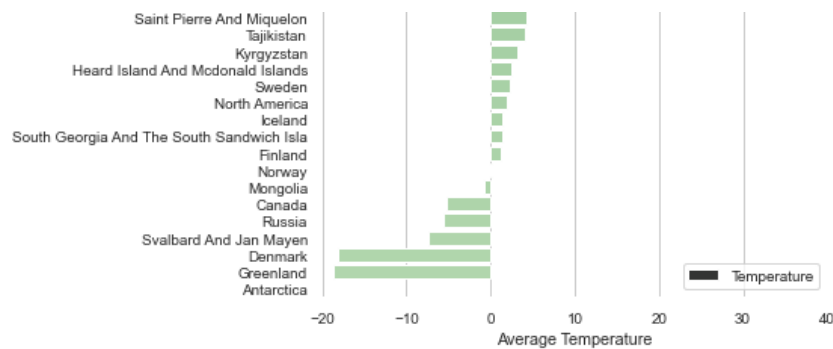
# Plot the temperature
sns.set_color_codes("pastel")
sns.barplot(x="AverageTemperature", y="Country", data=country_temp_asc,
            label="Temperature", palette="Greens_d")

# Informative axis label
ax.legend(ncol=2, loc="lower right", frameon=True)
ax.set(xlim=(-20, 40), ylabel="",
        xlabel="Average Temperature")
sns.despine(left=True, bottom=True)
```



Latin America		
Turks And Caicos Islands		
Cambodia		
Saint Lucia		
Solomon Islands		
American Samoa		
Saint Martin		
Sint Maarten		
Anguilla		
Saint Barthelemy		
Singapore		
Barbados		
Philippines		
Antigua And Barbuda		
Montserrat		
Mayotte		
British Virgin Islands		
Virgin Islands		
Samoa		
Guadeloupe		
Haiti		
Panama		
Yemen		
Suriname		
Trinidad And Tobago		
Martinique		
Dominica		
French Polynesia		
Jamaica		
Sierra Leone		
Thailand		
Timor Leste		
Nicaragua		
Côte D'Ivoire		
Guyana		
Bahrain		
Saint Kitts And Nevis		
Malaysia		
Sao Tome And Principe		
Christmas Island		
French Guiana		
Comoros		
Indonesia		
Costa Rica		
Dominican Republic		
Saudi Arabia		
Guinea		
Puerto Rico		
Cuba		
Liberia		
Baker Island		
Central African Republic		
Bahamas		
Kuwait		
Niue		
Belize		
Fiji		
Venezuela		
Equatorial Guinea		
El Salvador		
Bangladesh		
Colombia		
Brazil		
Honduras		
Congo		
Papua New Guinea		
Cameroon		
Gabon		
Cape Verde		
Kenya		
Africa		
India		
Congo (Democratic Republic Of The)		
Burma		
Vietnam		
Vietnam		
Mozambique		
Laos		
Mauritius		
Reunion		
Paraguay		
Tonga		
Palestina		
Guatemala		
Uganda		
Ethiopia		
Algeria		
Madagascar		
New Caledonia		
Hong Kong		
Egypt		
Macao		
Tanzania		
Western Sahara		
Libya		
Taiwan		
Botswana		
Ecuador		
Angola		
South America		





Теперь анализируем некоторые из самых жарких и прохладных стран и посмотрим, наблюдается ли повышение температуры с годами. Я хотел бы проанализировать Гренландию, Данию, Россию, Джибути и Мали

```
In [16]: GlobalTempCountry.rename(columns = {'dt':'Date'}, inplace = True)
```

```
In [17]: GlobalTempCountry['Year'] = pd.DatetimeIndex(GlobalTempCountry['Date']).year
GlobalTempCountry['Month'] = pd.DatetimeIndex(GlobalTempCountry['Date']).month
year_country = GlobalTempCountry.groupby(by = ['Year', 'Country']).mean().reset_index()
Russia = year_country.loc[year_country['Country'] == 'Russia',:]
Greenland = year_country.loc[year_country['Country'] == 'Greenland',:]
Denmark = year_country.loc[year_country['Country'] == 'Denmark',:]
Djibouti = year_country.loc[year_country['Country'] == 'Djibouti',:]
Mali = year_country.loc[year_country['Country'] == 'Mali',:]
Norway = year_country.loc[year_country['Country'] == 'Norway',:]
fig4 = go.Figure()
for template in ["plotly_dark"]:
    fig4.add_trace(go.Scatter(x=Russia['Year'], y=Russia['AverageTemperature'],
                             mode='lines',
                             name='Russia',
                             marker_color='#00CD66'))
    fig4.add_trace(go.Scatter(x=Greenland['Year'], y=Greenland['AverageTemperature'],
                             mode='lines',
                             name='Greenland',
                             marker_color='#FF4040'))
    fig4.add_trace(go.Scatter(x=Denmark['Year'], y=Denmark['AverageTemperature'],
                             mode='lines',
                             name='Denmark',
                             marker_color='#FFFF00'))
    fig4.add_trace(go.Scatter(x=Mali['Year'], y=Mali['AverageTemperature'],
                             mode='lines',
                             name='Mali',
                             marker_color='#EE82EE'))
    fig4.add_trace(go.Scatter(x=Djibouti['Year'], y=Djibouti['AverageTemperature'],
                             mode='lines',
                             name='Djibouti',
                             marker_color='#98F5FF'))
    fig4.add_trace(go.Scatter(x=Norway['Year'], y=Norway['AverageTemperature'],
                             mode='lines',
                             name='Norway',
                             marker_color='#E9967A'))
    fig4.update_layout(
        height=800,
        xaxis_title='Years',
        yaxis_title='Temperature in degree',
        title_text='Average Temperature Over the Years for the Following Countries',
        template=template)

fig4.show()
```



Из графика видно, что:

Гренландия имеет разницу в 3,77 градуса за данные годы.

Денамрк имеет разницу в 3,7 градуса за данные годы.

В России за эти годы разница составляет 4,78 градуса.

Норвегия имеет разницу в 5,2 градуса за данные годы.

В то время страны с высокими температурами:

Мали имеет разницу в 1,99 градуса за данные годы.

Джибути имеет разницу в 0,89 градуса за данные годы.

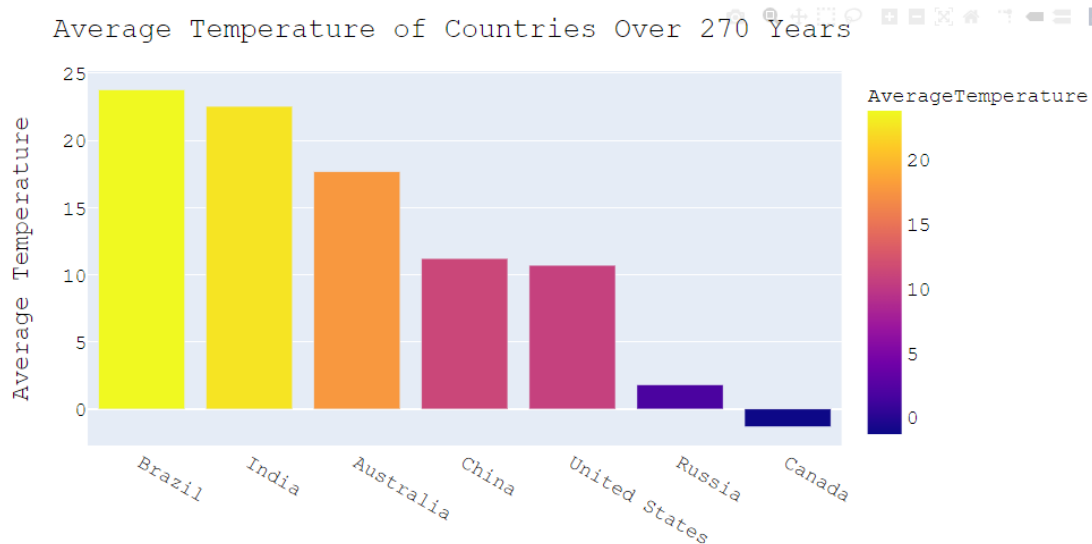
Таким образом, мы можем определенно сказать, что холодные места становятся все жарче.

Средняя температура данных 7 стран:

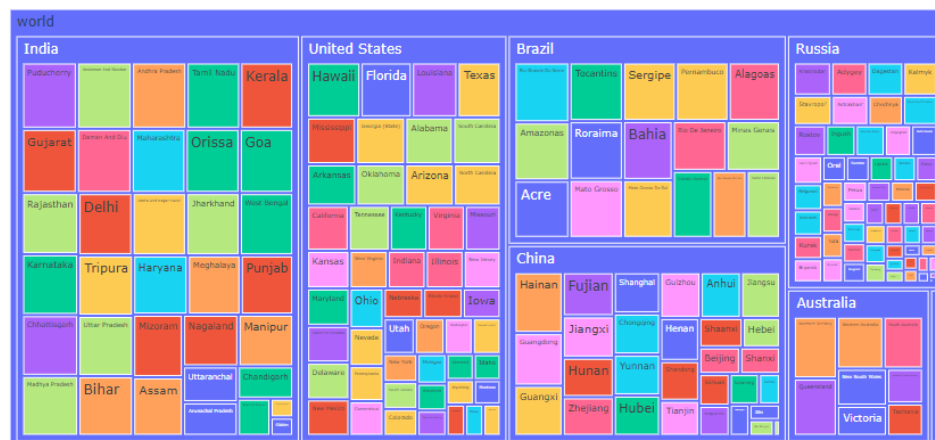
```
In [19]: country_temp_asc = GlobalTempState.groupby(by=['Country']).mean().reset_index().sort_values('AverageTemperature', ascending=False).reset_index()
country_temp_asc
plt.figure(figsize=(20,10))

fig5 = px.bar(country_temp_asc, x='Country', y='AverageTemperature', color='AverageTemperature')

fig5.update_layout(
    title='Average Temperature of Countries Over 270 Years',
    xaxis_title='Years',
    yaxis_title='Average Temperature',
    font=dict(
        family='Courier New',
        size=18,
        color='black'
    )
)
fig5.show()
```



```
In [20]: country_state_temp = GlobalTempState.groupby(by = ['Country', 'State']).mean().reset_index().sort_values('AverageTemperature', ascending=False).
country_state_temp
country_state_temp["world"] = "world" # in order to have a single root node
fig6 = px.treemap(country_state_temp.head(200), path=['world', 'Country', 'State'], values='AverageTemperature',
    color='State', color_continuous_scale='RdBu')
fig6.show()
```



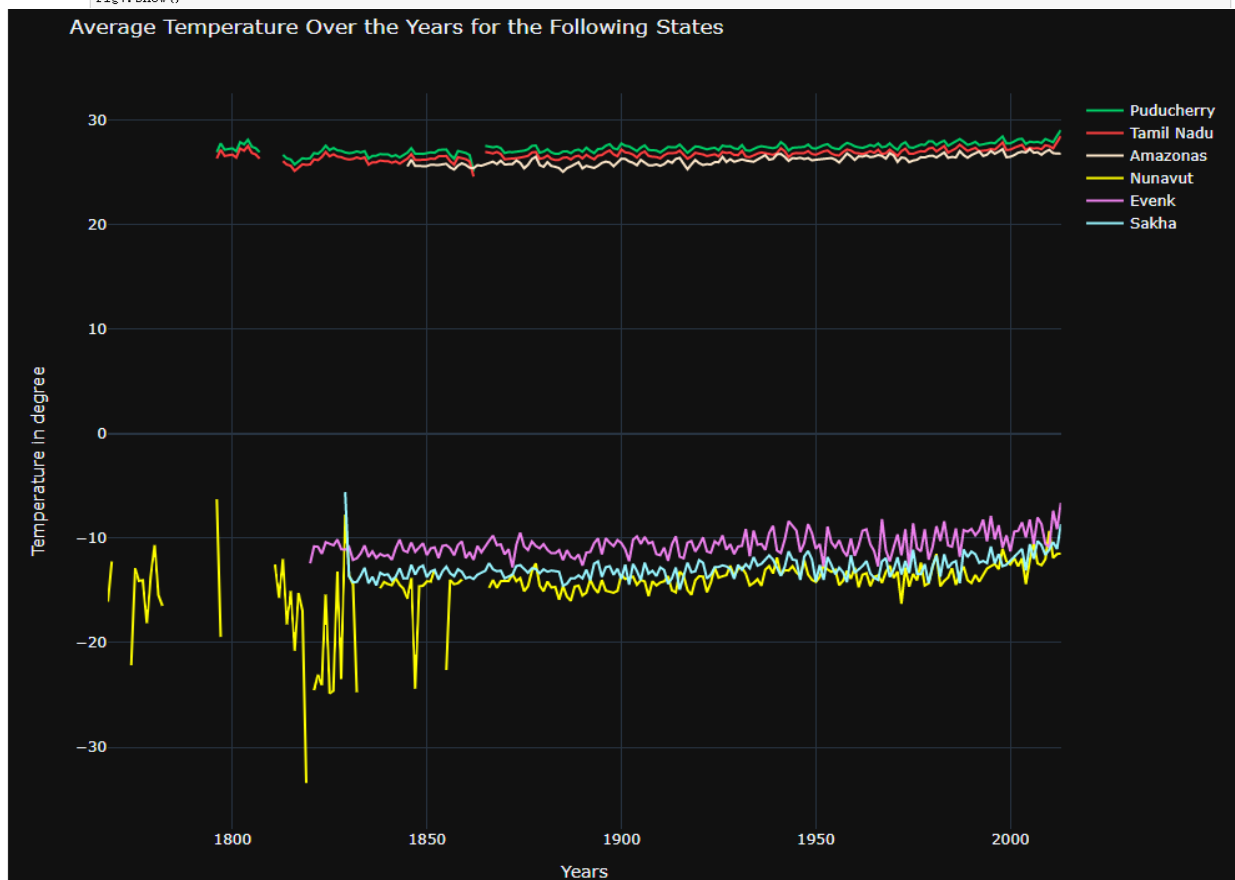
Теперь рассмотрим некоторые из самых горячих и холодных состояний, чтобы увидеть изменение температуры. Рассмотрим Пудучерри, Тамилнад из Индии, Амазонас из Бразилии, Нунавут из Канады, эвенкию и Саху из России.

```

In [21]: GlobalTempState['Year'] = pd.DatetimeIndex(GlobalTempState['Date']).year
year_state = GlobalTempState.groupby(by = ['Year', 'State']).mean().reset_index()
Puducherry = year_state.loc[year_state['State'] == 'Puducherry',:]
Tamil_Nadu = year_state.loc[year_state['State'] == 'Tamil Nadu',:]
Amazonas = year_state.loc[year_state['State'] == 'Amazonas',:]
Munavut = year_state.loc[year_state['State'] == 'Munavut',:]
Evenk = year_state.loc[year_state['State'] == 'Evenk',:]
Sakha = year_state.loc[year_state['State'] == 'Sakha',:]
fig7 = go.Figure()
for template in ["plotly_dark"]:
    fig7.add_trace(go.Scatter(x=Puducherry['Year'], y=Puducherry['AverageTemperature'],
                             mode='lines',
                             name='Puducherry',
                             marker_color='#00CD66'))
    fig7.add_trace(go.Scatter(x=Tamil_Nadu['Year'], y=Tamil_Nadu['AverageTemperature'],
                             mode='lines',
                             name='Tamil Nadu',
                             marker_color='#FF4040'))
    fig7.add_trace(go.Scatter(x=Amazonas['Year'], y=Amazonas['AverageTemperature'],
                             mode='lines',
                             name='Amazonas',
                             marker_color='#FCEB6C'))
    fig7.add_trace(go.Scatter(x=Munavut['Year'], y=Munavut['AverageTemperature'],
                             mode='lines',
                             name='Munavut',
                             marker_color='#FFFF00'))
    fig7.add_trace(go.Scatter(x=Evenk['Year'], y=Evenk['AverageTemperature'],
                             mode='lines',
                             name='Evenk',
                             marker_color='#EB82EE'))
    fig7.add_trace(go.Scatter(x=Sakha['Year'], y=Sakha['AverageTemperature'],
                             mode='lines',
                             name='Sakha',
                             marker_color='#98F5FF'))
    fig7.update_layout(
        height=800,
        xaxis_title='Years',
        yaxis_title='Temperature in degree',
        title_text='Average Temperature Over the Years for the Following States',
        template=template)

fig7.show()

```



Видно, что:
Пудучерри имеет разницу в 1,88 градуса за эти годы.

Тамилнад имеет разницу в 1,56 градуса за эти годы.

Амазонас имеет разницу в 1,28 градуса за эти годы.

Нунавут имеет разницу в 4,6 градуса за эти годы.

Эвенкия имеет разницу в 5,8 градуса за эти годы.

Саха имеет разницу в 3,05 градуса за эти годы.

Можем сказать, что в холодных странах температура поднимается быстрее, чем в жарких, который считается признаком глобального потепления.

Список литературы

[1] Гапанюк Ю. Е. Лабораторная работа «Разведочный анализ данных. Исследование и визуализация данных» [Электронный ресурс] // GitHub. — 2019. — Режим доступа: https://github.com/ugapanyuk/ml_course/wiki/LAB_EDA_VISUALIZATION (дата обращения: 13.02.2019)

[2] <https://www.kaggle.com/datasets>