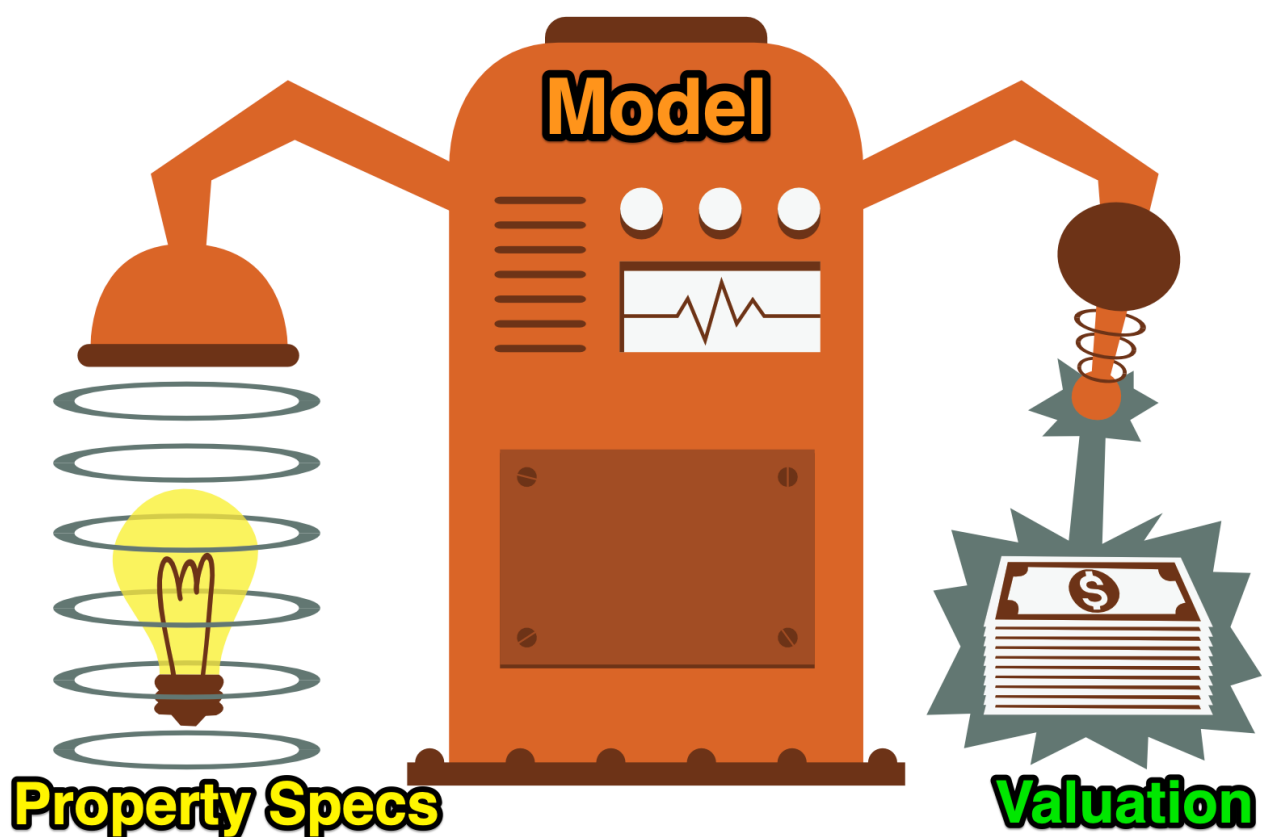# Setup and Context

## Introduction

Welcome to Boston Massachusetts in the 1970s! Imagine you're working for a real estate development company. Your company wants to value any residential project before they start. You are tasked with building a model that can provide a price estimate based on a home's characteristics like:

- The number of rooms
- The distance to employment centres
- How rich or poor the area is
- How many students there are per teacher in local schools etc

**To accomplish your task you will:**

1. Analyse and explore the Boston house price data
2. Split your data for training and testing
3. Run a Multivariable Regression
4. Evaluate how your model's coefficients and residuals
5. Use data transformation to improve your model performance
6. Use your model to estimate a property price

## Upgrade plotly (only Google Colab Notebook)

Google Colab may not be running the latest version of plotly. If you're working in Google Colab, uncomment the line below, run the cell, and restart your notebook server.

In [30]:

```python
# %pip install --upgrade plotly
```

## Import Statements

In [31]:

```python
import pandas as pd
import numpy as np

import seaborn as sns
import plotly.express as px
import matplotlib.pyplot as plt

from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
# TODO: Add missing import statements
```

## Notebook Presentation

In [32]:

```python
pd.options.display.float_format = '{:,.2f}'.format
```

# Load the Data

The first column in the .csv file just has the row numbers, so it will be used as the index.

In [33]:

```python
data = pd.read_csv('boston.csv', index_col=0)
```

## Understand the Boston House Price Dataset

**Characteristics:**

```
:Number of Instances: 506

:Number of Attributes: 13 numeric/categorical predictive. The Median Value (attrib
ute 14) is the target.

:Attribute Information (in order):
    1. CRIM      per capita crime rate by town
```

```
     2.  ZN        proportion of residential land zoned for lots over 25,000 sq.ft.
     3.  INDUS     proportion of non-retail business acres per town
     4.  CHAS      Charles River dummy variable (= 1 if tract bounds river; 0 otherwi
   se)
     5.  NOX       nitric oxides concentration (parts per 10 million)
     6.  RM        average number of rooms per dwelling
     7.  AGE       proportion of owner-occupied units built prior to 1940
     8.  DIS       weighted distances to five Boston employment centres
     9.  RAD       index of accessibility to radial highways
     10. TAX        full-value property-tax rate per $10,000
     11. PTRATIO   pupil-teacher ratio by town
     12. B          1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
     13. LSTAT     % lower status of the population
     14. PRICE      Median value of owner-occupied homes in $1000's

   :Missing Attribute Values: None

   :Creator: Harrison, D. and Rubinfeld, D.L.
```

**This is a copy of [UCI ML housing dataset](). This dataset was taken from the StatLib library which is maintained at Carnegie Mellon University. You can find the [original research paper here]().**

# Preliminary Data Exploration 🏠

**Challenge**

- **What is the shape of `data`?**
- **How many rows and columns does it have?**
- **What are the column names?**
- **Are there any NaN values or duplicates?**

In [34]:

```
data.shape
```

Out[34]:

```
(506, 14)
```

In [35]:

```
data.columns
```

Out[35]:

```
Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
       'PTRATIO', 'B', 'LSTAT', 'PRICE'],
      dtype='object')
```

In [36]:

```
data.dtypes
```

Out[36]:

| | 0 |
| --- | --- |
| CRIM | float64 |
| ZN | float64 |
| INDUS | float64 |
| CHAS | float64 |
| NOX | float64 |

| | |
|---:|:---|
| **RM** | float64 |
| **AGE** | float64 |
| **DIS** | float64 |
| **RAD** | float64 |
| **TAX** | float64 |
| **PTRATIO** | float64 |
| **B** | float64 |
| **LSTAT** | float64 |
| **PRICE** | float64 |

**dtype: object**

## Data Cleaning - Check for Missing Values and Duplicates

In [37]:

```
data.isna().sum()
```

Out[37]:

| | 0 |
|---:|:---|
| **CRIM** | 0 |
| **ZN** | 0 |
| **INDUS** | 0 |
| **CHAS** | 0 |
| **NOX** | 0 |
| **RM** | 0 |
| **AGE** | 0 |
| **DIS** | 0 |
| **RAD** | 0 |
| **TAX** | 0 |
| **PTRATIO** | 0 |
| **B** | 0 |
| **LSTAT** | 0 |
| **PRICE** | 0 |

**dtype: int64**

In [38]:

```
data.duplicated().sum()
```

Out[38]:

0

In [39]:

```
data.sample(10)
```

Out[39]:

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | PRICE |
|---:|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **238** | 0.08 | 30.00 | 4.93 | 0.00 | 0.43 | 6.48 | 18.50 | 6.19 | 6.00 | 300.00 | 16.60 | 379.41 | 6.36 | 23.70 |

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | PRICE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 459 | 6.80 | 0.00 | 18.10 | 0.00 | 0.71 | 6.08 | 84.40 | 2.72 | 24.00 | 666.00 | 20.20 | 396.90 | 14.70 | 20.00 |
| 419 | 11.81 | 0.00 | 18.10 | 0.00 | 0.72 | 6.82 | 76.50 | 1.79 | 24.00 | 666.00 | 20.20 | 48.45 | 22.74 | 8.40 |
| 498 | 0.24 | 0.00 | 9.69 | 0.00 | 0.58 | 6.02 | 65.30 | 2.41 | 6.00 | 391.00 | 19.20 | 396.90 | 12.92 | 21.20 |
| 87 | 0.07 | 0.00 | 4.49 | 0.00 | 0.45 | 6.12 | 56.80 | 3.75 | 3.00 | 247.00 | 18.50 | 395.15 | 8.44 | 22.20 |
| 225 | 0.53 | 0.00 | 6.20 | 0.00 | 0.50 | 8.72 | 83.00 | 2.89 | 8.00 | 307.00 | 17.40 | 382.00 | 4.63 | 50.00 |
| 67 | 0.06 | 12.50 | 6.07 | 0.00 | 0.41 | 5.88 | 21.40 | 6.50 | 4.00 | 345.00 | 18.90 | 396.21 | 8.10 | 22.00 |
| 405 | 67.92 | 0.00 | 18.10 | 0.00 | 0.69 | 5.68 | 100.00 | 1.43 | 24.00 | 666.00 | 20.20 | 384.97 | 22.98 | 5.00 |
| 150 | 1.66 | 0.00 | 19.58 | 0.00 | 0.87 | 6.12 | 97.30 | 1.62 | 5.00 | 403.00 | 14.70 | 372.80 | 14.10 | 21.50 |
| 249 | 0.19 | 22.00 | 5.86 | 0.00 | 0.43 | 6.72 | 17.50 | 7.83 | 7.00 | 330.00 | 19.10 | 393.74 | 6.56 | 26.20 |

# Descriptive Statistics

**Challenge**

- **How many students are there per teacher on average?**
- **What is the average price of a home in the dataset?**
- **What is the `CHAS` feature?**
- **What are the minimum and the maximum value of the `CHAS` and why?**
- **What is the maximum and the minimum number of rooms per dwelling in the dataset?**

In [40]:

```
data.describe()
```

Out[40]:

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | PRICE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 | 506.00 |
| mean | 3.61 | 11.36 | 11.14 | 0.07 | 0.55 | 6.28 | 68.57 | 3.80 | 9.55 | 408.24 | 18.46 | 356.67 | 12.65 | 22.53 |
| std | 8.60 | 23.32 | 6.86 | 0.25 | 0.12 | 0.70 | 28.15 | 2.11 | 8.71 | 168.54 | 2.16 | 91.29 | 7.14 | 9.20 |
| min | 0.01 | 0.00 | 0.46 | 0.00 | 0.39 | 3.56 | 2.90 | 1.13 | 1.00 | 187.00 | 12.60 | 0.32 | 1.73 | 5.00 |
| 25% | 0.08 | 0.00 | 5.19 | 0.00 | 0.45 | 5.89 | 45.02 | 2.10 | 4.00 | 279.00 | 17.40 | 375.38 | 6.95 | 17.02 |
| 50% | 0.26 | 0.00 | 9.69 | 0.00 | 0.54 | 6.21 | 77.50 | 3.21 | 5.00 | 330.00 | 19.05 | 391.44 | 11.36 | 21.20 |
| 75% | 3.68 | 12.50 | 18.10 | 0.00 | 0.62 | 6.62 | 94.07 | 5.19 | 24.00 | 666.00 | 20.20 | 396.23 | 16.96 | 25.00 |
| max | 88.98 | 100.00 | 27.74 | 1.00 | 0.87 | 8.78 | 100.00 | 12.13 | 24.00 | 711.00 | 22.00 | 396.90 | 37.97 | 50.00 |

In [40]:

In [41]:

```
data.sample(10)
```

Out[41]:

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | PRICE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 9 | 0.17 | 12.50 | 7.87 | 0.00 | 0.52 | 6.00 | 85.90 | 6.59 | 5.00 | 311.00 | 15.20 | 386.71 | 17.10 | 18.90 |
| 489 | 0.18 | 0.00 | 27.74 | 0.00 | 0.61 | 5.41 | 98.30 | 1.76 | 4.00 | 711.00 | 20.10 | 344.05 | 23.97 | 7.00 |
| 405 | 67.92 | 0.00 | 18.10 | 0.00 | 0.69 | 5.68 | 100.00 | 1.43 | 24.00 | 666.00 | 20.20 | 384.97 | 22.98 | 5.00 |
| 454 | 9.51 | 0.00 | 18.10 | 0.00 | 0.71 | 6.73 | 94.10 | 2.50 | 24.00 | 666.00 | 20.20 | 6.68 | 18.71 | 14.90 |
| 301 | 0.04 | 34.00 | 6.09 | 0.00 | 0.43 | 6.59 | 40.40 | 5.49 | 7.00 | 329.00 | 16.10 | 395.75 | 9.50 | 22.00 |

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | PRICE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 257 | 0.61 | 20.00 | 3.97 | 0.00 | 0.65 | 8.70 | 86.90 | 1.80 | 5.00 | 264.00 | 13.00 | 389.70 | 5.12 | 50.00 |
| ~~317~~ | ~~0.25~~ | ~~0.00~~ | ~~9.90~~ | ~~0.00~~ | ~~0.54~~ | ~~5.78~~ | ~~71.70~~ | ~~4.03~~ | ~~4.00~~ | ~~304.00~~ | ~~18.40~~ | ~~396.90~~ | ~~15.94~~ | ~~19.80~~ |
| 116 | 0.13 | 0.00 | 10.01 | 0.00 | 0.55 | 6.18 | 72.50 | 2.73 | 6.00 | 432.00 | 17.80 | 393.30 | 12.04 | 21.20 |
| 179 | 0.06 | 0.00 | 2.46 | 0.00 | 0.49 | 6.98 | 58.40 | 2.83 | 3.00 | 193.00 | 17.80 | 396.90 | 5.04 | 37.20 |
| 286 | 0.02 | 80.00 | 1.76 | 0.00 | 0.39 | 6.23 | 31.50 | 9.09 | 1.00 | 241.00 | 18.20 | 341.60 | 12.93 | 20.10 |

# Visualise the Features

**Challenge:** Having looked at some descriptive statistics, visualise the data for your model. Use  Seaborn's  `.displot()`  to create a bar chart and superimpose the Kernel Density Estimate (KDE) for the following variables:

- **PRICE: The home price in thousands.**
- **RM: the average number of rooms per owner unit.**
- **DIS: the weighted distance to the 5 Boston employment centres i.e., the estimated length of the commute.**
- **RAD: the index of accessibility to highways.**

Try setting the  `aspect`  parameter to  `2`  for a better picture.

**What do you notice in the distributions of the data?**

**House Prices 🏠**

In [42]:

```
sns.displot(data['PRICE'],
            bins=50,
            aspect=2,
            kde=True,
            color='#2196f3')

plt.title(f'1970s Home Values in Boston. Average: ${(1000*data.PRICE.mean()):.6}')
plt.xlabel('Price in 000s')
plt.ylabel('Nr. of Homes')

plt.show()
```



**Distance to Employment - Length of Commute 🚗**

```
sns.displot(data['DIS'],
            bins=50,
            aspect=2,
            kde=True,
            color='#2196f3')
plt.title("Distance to the employment centers")
plt.xlabel("Distance to employment")
plt.ylabel("Nr. of Houses")
plt.show()
```
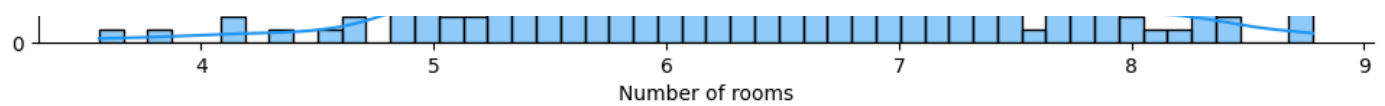


Distance to the employment centers

## Number of Rooms

```
sns.displot(data['RM'],
            bins=50,
            aspect=2,
            kde=True,
            color='#2196f3')
plt.title("Number of rooms per employment center")
plt.xlabel("Number of rooms")
plt.ylabel("Nr. of Houses")
plt.show()
```



Number of rooms per employment center

Number of rooms

**Access to Highways** 🔗

```python
sns.displot(data['RAD'],
            bins=50,
            aspect=2,
            kde=True,
            color='#2196f3')

plt.title('Accesibility from homes to highways')
plt.xlabel('Accesibility')
plt.ylabel('Nr. houses')
```

Text(5.069444444444445, 0.5, 'Nr. houses')



**Next to the River?** 🔗

**Challenge**

Create a bar chart with plotly for CHAS to show many more homes are away from the river versus next to it. The bar chart should look something like this:

Next to Charles River?



Any color will do

100

0

No          Yes

Property Located Next to the River?

**You can make your life easier by providing a list of values for the x-axis (e.g.,** `x=['No', 'Yes']` **)**

In [46]:

```
px.bar(x=['No', 'Yes'],
       y=data['CHAS'].value_counts(),
       labels={'x': 'Next to River?', 'y': 'Nr. of Homes'},)
```

In [46]:

# Understand the Relationships in the Data

### Run a Pair Plot

**Challenge**

There might be some relationships in the data that we should know about. Before you run the code, make some predictions:

- What would you expect the relationship to be between pollution (NOX) and the distance to employment (DIS)?
- What kind of relationship do you expect between the number of rooms (RM) and the home value (PRICE)?
- What about the amount of poverty in an area (LSTAT) and home prices?

Run a **Seaborn** `.pairplot()` to visualise all the relationships at the same time. Note, this is a big task and can take 1-2 minutes! After it's finished check your intuition regarding the questions above on the `pairplot`.

In [47]:
```
sns.pairplot(data)
```

Output hidden; open in https://colab.research.google.com to view.

In [47]:

**Challenge**

Use **Seaborn's** `.jointplot()` to look at some of the relationships in more detail. Create a jointplot for:

- DIS and NOX
- INDUS vs NOX
- LSTAT vs RM
- LSTAT vs PRICE
- RM vs PRICE

Try adding some opacity or `alpha` to the scatter plots using keyword arguments under `joint_kws`.

### Distance from Employment vs. Pollution

**Challenge:**

Compare DIS (Distance from employment) with NOX (Nitric Oxide Pollution) using Seaborn's `.jointplot()`. Does pollution go up or down as the distance increases?

In [48]:
```
sns.jointplot(x=data['DIS'],
              y=data['NOX'],)
```

Out[48]:

```
<seaborn.axisgrid.JointGrid at 0x7d305a991f90>
```

**Proportion of Non-Retail Industry 🏭🏭🏭 versus Pollution**

**Challenge:**

**Compare INDUS (the proportion of non-retail industry i.e., factories) with NOX (Nitric Oxide Pollution) using Seaborn's** `.jointplot()` **. Does pollution go up or down as there is a higher proportion of industry?**
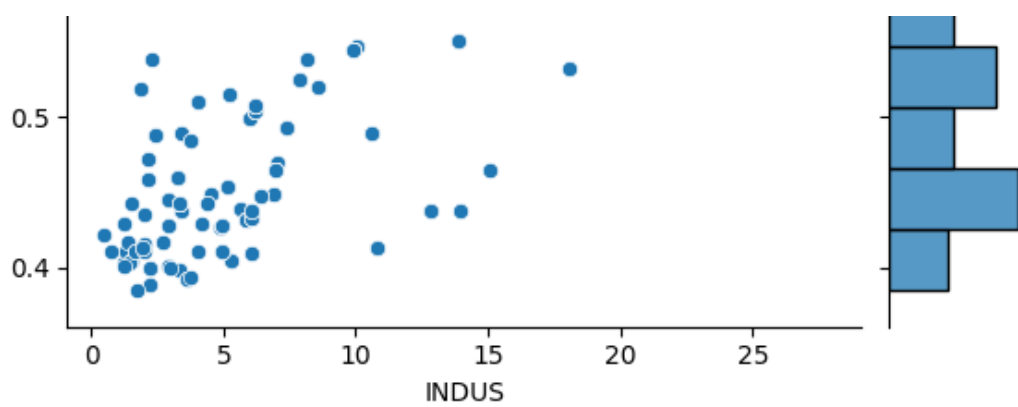
In [49]:

```
sns.jointplot(x=data['INDUS'],
              y=data['NOX'],)
```

Out[49]:

```
<seaborn.axisgrid.JointGrid at 0x7d30525c65c0>
```
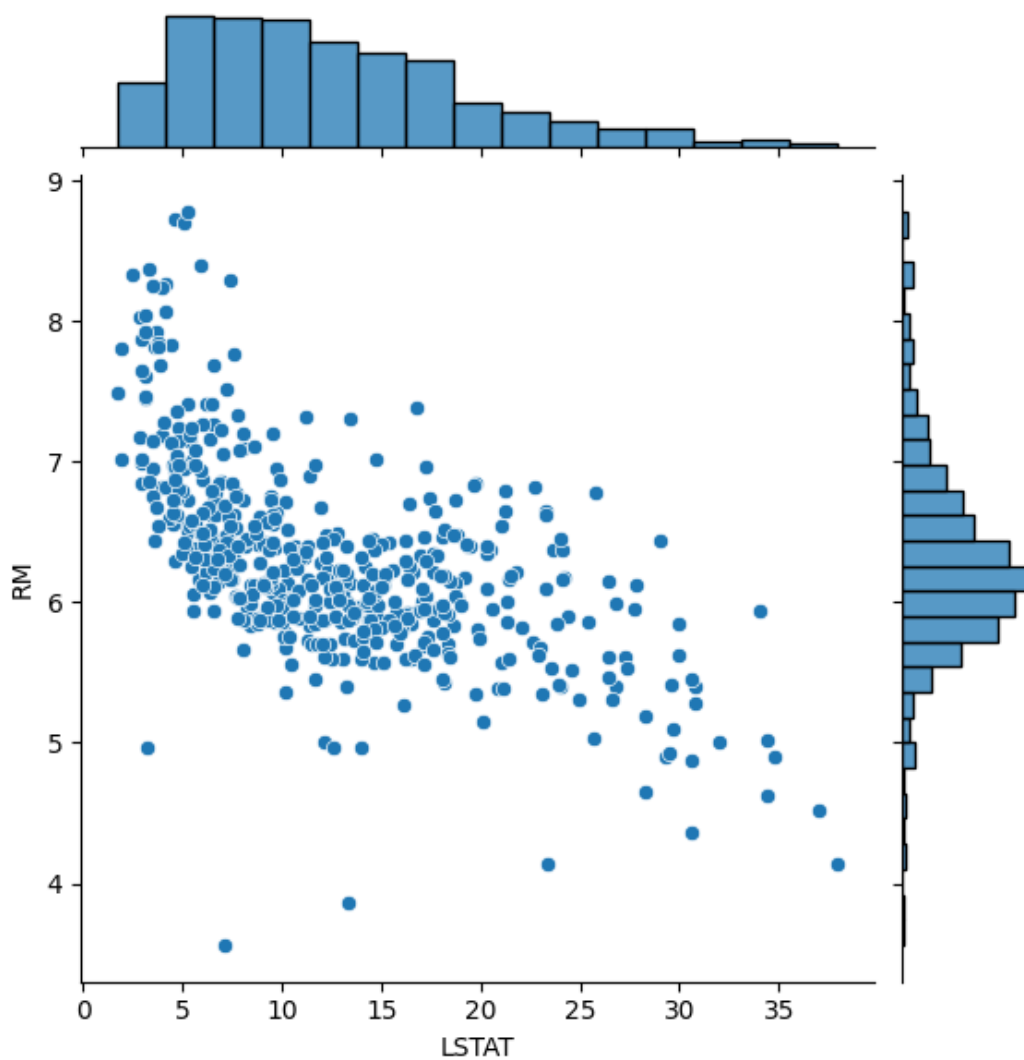
**% of Lower Income Population vs Average Number of Rooms**

**Challenge**

Compare LSTAT (proportion of lower-income population) with RM (number of rooms) using Seaborn's `.jointplot()` . How does the number of rooms per dwelling vary with the poverty of area? Do homes have more or fewer rooms when LSTAT is low?

In [50]:

```
sns.jointplot(x=data['LSTAT'],
              y=data['RM'])
```

Out[50]:

```
<seaborn.axisgrid.JointGrid at 0x7d30525f3d60>
```



**% of Lower Income Population versus Home Price**

### Challenge

Compare LSTAT with PRICE using Seaborn's `.jointplot()`. How does the proportion of the lower-income population in an area affect home prices?

In [51]:

```
sns.jointplot(x=data['LSTAT'],
              y=data['PRICE'])
```

Out[51]:

```
<seaborn.axisgrid.JointGrid at 0x7d30525f1990>
```



### Number of Rooms versus Home Value

### Challenge

Compare RM (number of rooms) with PRICE using Seaborn's `.jointplot()`. You can probably guess how the number of rooms affects home prices. ☺

In [52]:

```
sns.jointplot(x=data['RM'],
              y=data['PRICE'])
```

Out[52]:

```
<seaborn.axisgrid.JointGrid at 0x7d304f53be50>
```

## Split Training & Test Dataset

We *can't* use all 506 entries in our dataset to train our model. The reason is that we want to evaluate our model on data that it hasn't seen yet (i.e., out-of-sample data). That way we can get a better idea of its performance in the real world.

**Challenge**

- **Import the `train_test_split()` function from sklearn**
- **Create 4 subsets: X_train, X_test, y_train, y_test**
- **Split the training and testing data roughly 80/20.**
- **To get the same random split every time you run your notebook use `random_state=10`. This helps us get the same results every time and avoid confusion while we're learning.**

**Hint: Remember, your target is your home PRICE, and your features are all the other columns you'll use to predict the price.**

In [53]:

```
target =data['PRICE']
features = data.drop(['PRICE'], axis=1)

X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2, ran
dom_state=10)
```

In [54]:

```
train_pct = 100*len(X_train)/len(features)
print(f'Training data is {train_pct:.3}% of the total data.')

test_pct = 100*X_test.shape[0]/features.shape[0]
print(f'Test data makes up the remaining {test_pct:0.3}%.')
```

```
Training data is 79.8% of the total data.
Test data makes up the remaining 20.2%.
```

# Multivariable Regression

In a previous lesson, we had a linear model with only a single feature (our movie budgets). This time we have a total of 13 features. Therefore, our Linear Regression model will have the following form:

$$PR\hat{I}CE = \theta_0 + \theta_1\,RM$$
$$+\,\theta_2\,NOX + \theta_3\,DIS$$
$$+\,\theta_4\,CHAS...$$
$$+\theta_{13}\,LSTAT$$

## Run Your First Regression

### Challenge

Use sklearn to run the regression on the training dataset. How high is the r-squared for the regression on the training data?

In [59]:

```
regr = LinearRegression()
regr.fit(X_train, y_train)
rsquared = regr.score(X_train, y_train)

print(f'Training data r-squared: {rsquared:.2}')
```

Training data r-squared: 0.75

In [54]:

## Evaluate the Coefficients of the Model

Here we do a sense check on our regression coefficients. The first thing to look for is if the coefficients have the expected sign (positive or negative).

**Challenge** Print out the coefficients (the thetas in the equation above) for the features. Hint: You'll see a nice table if you stick the coefficients in a DataFrame.

- We already saw that RM on its own had a positive relation to PRICE based on the scatter plot. Is RM's coefficient also positive?
- What is the sign on the LSAT coefficient? Does it match your intuition and the scatter plot above?
- Check the other coefficients. Do they have the expected sign?
- Based on the coefficients, how much more expensive is a room with 6 rooms compared to a room with 5 rooms? According to the model, what is the premium you would have to pay for an extra room?

In [60]:

```
regr_coef = pd.DataFrame(data=regr.coef_, index=X_train.columns, columns=['Coefficient']
)
regr_coef
```

Out[60]:

|        | Coefficient |
|--------|-------------|
| CRIM   | -0.13       |
| ZN     | 0.06        |
| INDUS  | -0.01       |
| CHAS   | 1.97        |
| NOX    | -16.27      |
| RM     | 3.11        |

| | Coefficient |
|---|---|
| AGE | 0.02 |
| DIS | -1.48 |
| RAD | 0.30 |
| TAX | -0.01 |
| PTRATIO | -0.82 |
| B | 0.01 |
| LSTAT | -0.58 |

In [61]:

```
premium = regr_coef.loc['RM'].values[0] * 1000   # i.e., ~3.11 * 1000
print(f'The price premium for having an extra room is ${premium:.5}')
```
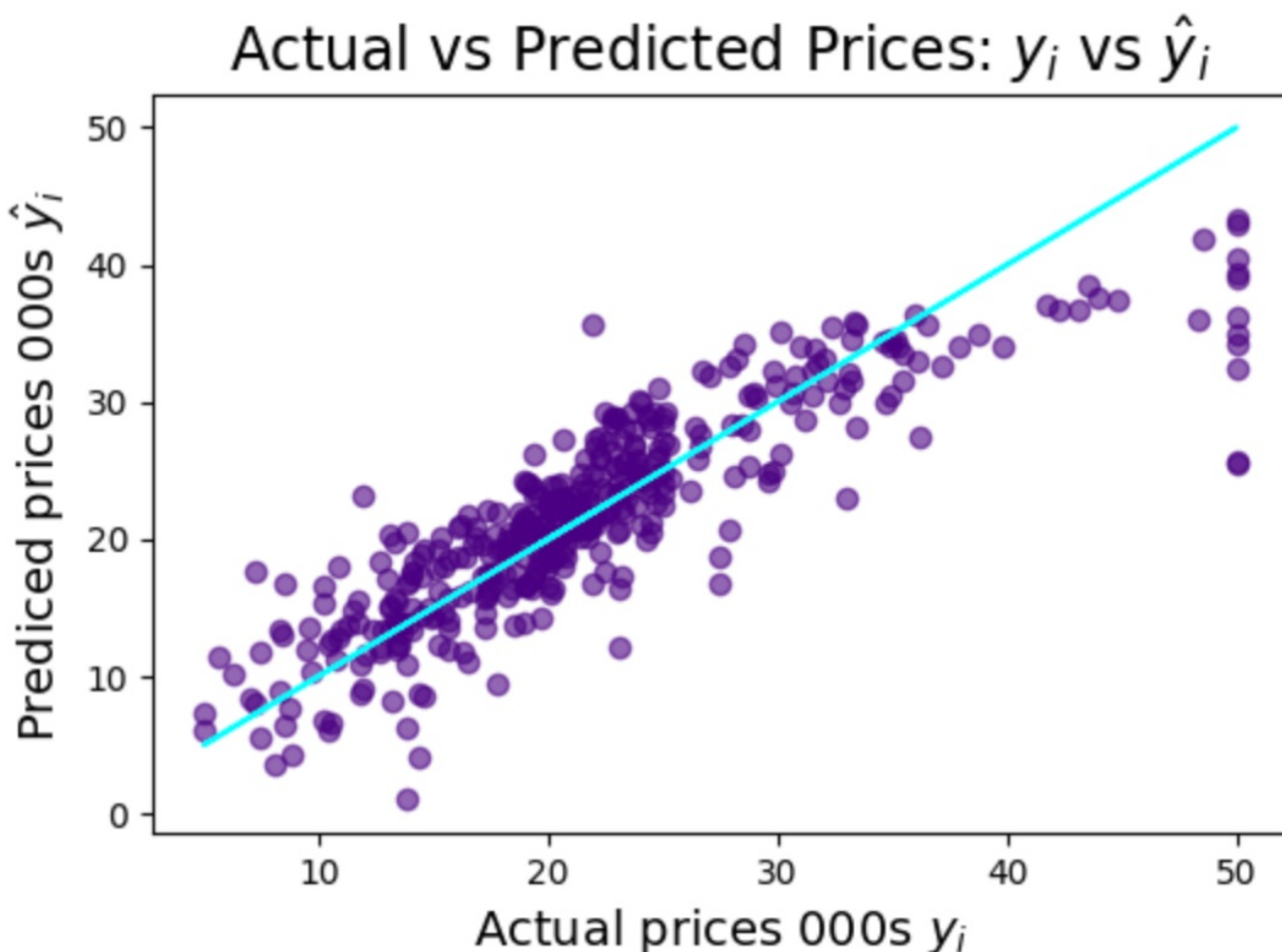
The price premium for having an extra room is $3108.5

## Analyse the Estimated Values & Regression Residuals

The next step is to evaluate our regression. How good our regression is depends not only on the r-squared. It also depends on the **residuals** - the difference between the model's predictions ( $\hat{y}_i$ ) and the true values ( $y_i$ ) inside `y_train` .

```
    predicted_values = regr.predict(X_train)
    residuals = (y_train - predicted_values)
```

**Challenge:** Create two scatter plots.

The first plot should be actual values ( `y_train` ) against the predicted value values:



The cyan line in the middle shows `y_train` against `y_train` . If the predictions had been 100% accurate then all the dots would be on this line. The further away the dots are from the line, the worse the prediction was. That makes the distance to the cyan line, you guessed it, our residuals ☺

The second plot should be the residuals against the predicted prices. Here's what we're looking for:
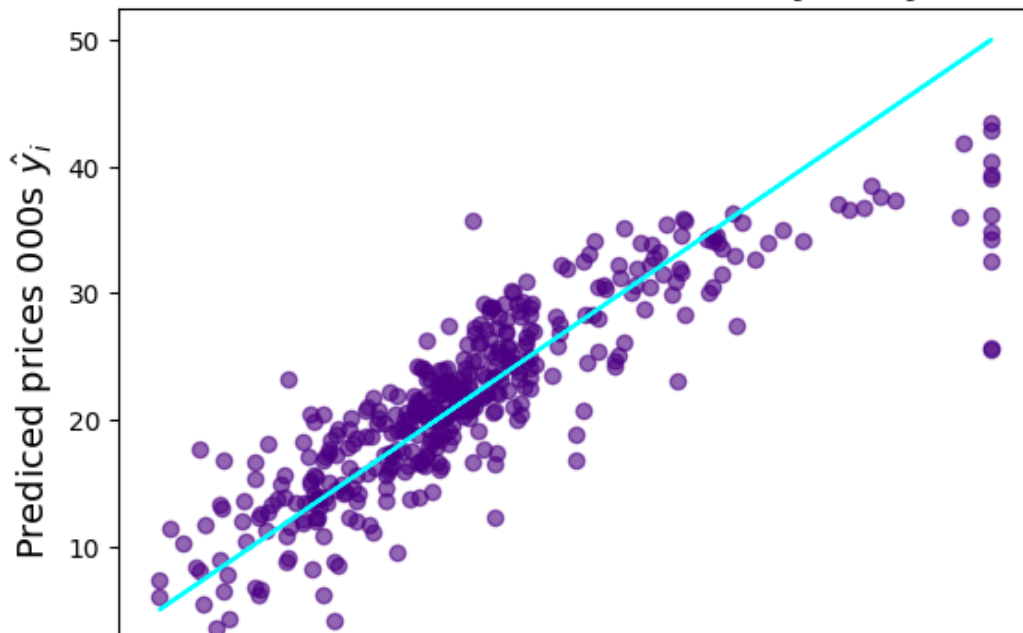
# Residuals vs Predicted Values

```python
predicted_vals = regr.predict(X_train)
residuals = (y_train - predicted_vals)
```
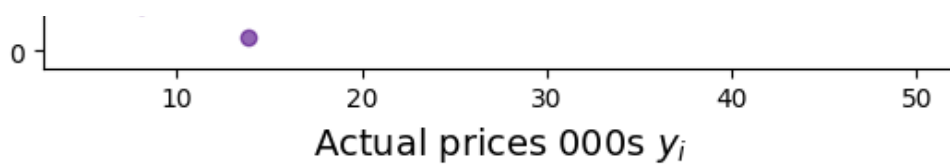
```python
# Original Regression of Actual vs. Predicted Prices
plt.figure(dpi=100)
plt.scatter(x=y_train, y=predicted_vals, c='indigo', alpha=0.6)
plt.plot(y_train, y_train, color='cyan')
plt.title(f'Actual vs Predicted Prices: $y _i$ vs $\hat y_i$', fontsize=17)
plt.xlabel('Actual prices 000s $y _i$', fontsize=14)
plt.ylabel('Prediced prices 000s $\hat y _i$', fontsize=14)
plt.show()
```

## Actual vs Predicted Prices: $y_i$ vs $\hat{y}_i$

Actual prices 000s $y_i$

```python
# Residuals vs Predicted values
plt.figure(dpi=100)
plt.scatter(x=predicted_vals, y=residuals, c='indigo', alpha=0.6)
plt.title('Residuals vs Predicted Values', fontsize=17)
plt.xlabel('Predicted Prices $\hat y _i$', fontsize=14)
plt.ylabel('Residuals', fontsize=14)
plt.show()
```
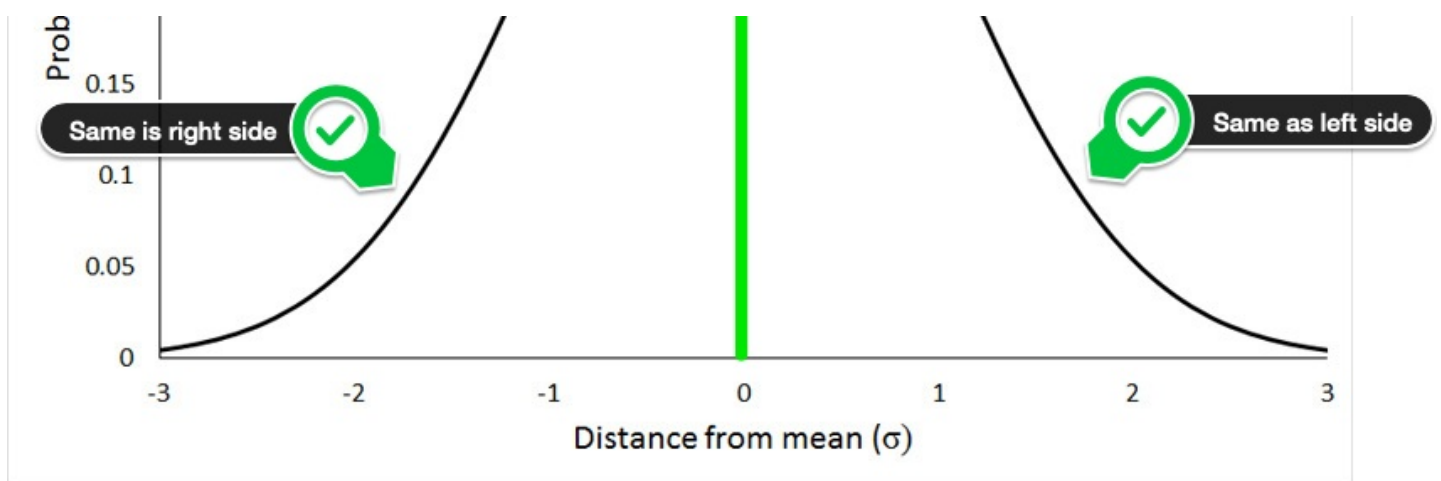


**Why do we want to look at the residuals? We want to check that they look random. Why? The residuals represent the errors of our model. If there's a pattern in our errors, then our model has a systematic bias.**

**We can analyse the distribution of the residuals. In particular, we're interested in the skew and the mean.**

**In an ideal case, what we want is something close to a normal distribution. A normal distribution has a skewness of 0 and a mean of 0. A skew of 0 means that the distribution is symmetrical - the bell curve is not lopsided or biased to one side. Here's what a normal distribution looks like:**

**Same is right side** ✓

✓ **Same as left side**

**Prob** — y-axis: 0, 0.05, 0.1, 0.15

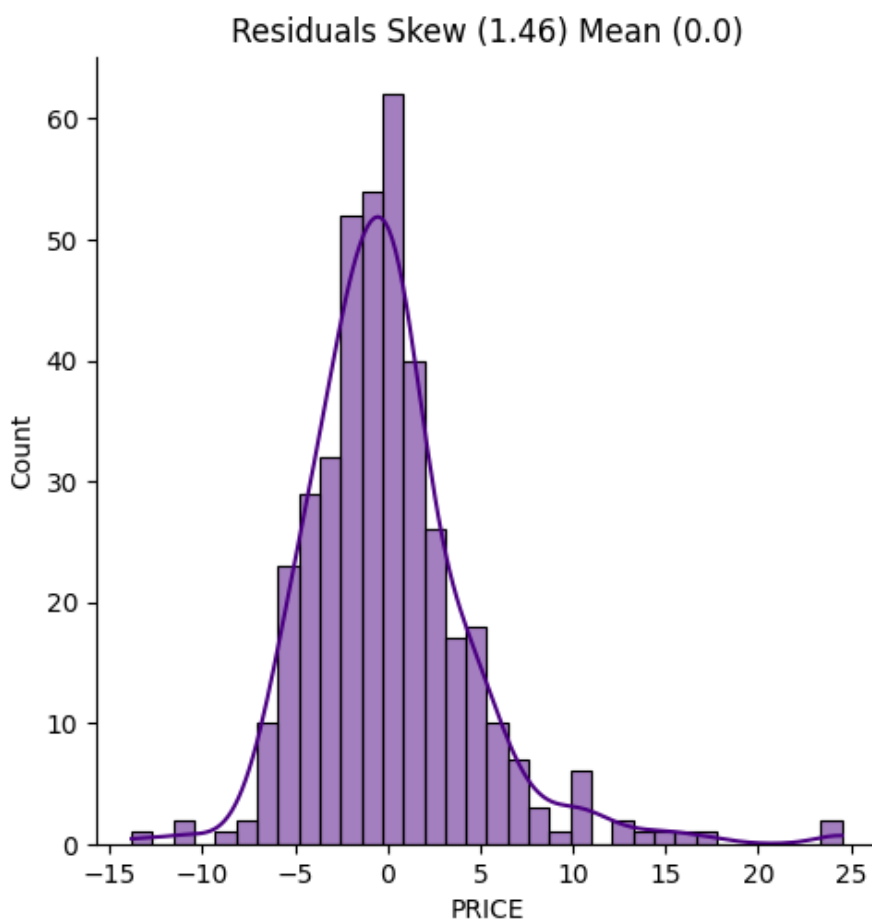x-axis: -3, -2, -1, 0, 1, 2, 3

**Distance from mean (σ)**

### Challenge

- Calculate the mean and the skewness of the residuals.
- Again, use Seaborn's `.displot()` to create a histogram and superimpose the Kernel Density Estimate (KDE)
- Is the skewness different from zero? If so, by how much?
- Is the mean different from zero?

In [65]:

```python
# Residual Distribution Chart
resid_mean = round(residuals.mean(), 2)
resid_skew = round(residuals.skew(), 2)

sns.displot(residuals, kde=True, color='indigo')
plt.title(f'Residuals Skew ({resid_skew}) Mean ({resid_mean})')
plt.show()
```



Residuals Skew (1.46) Mean (0.0)

In [54]:

# Data Transformations for a Better Fit

**We have two options at this point:**

1. **Change our model entirely. Perhaps a linear model is not appropriate.**
2. **Transform our data to make it fit better with our linear model.**

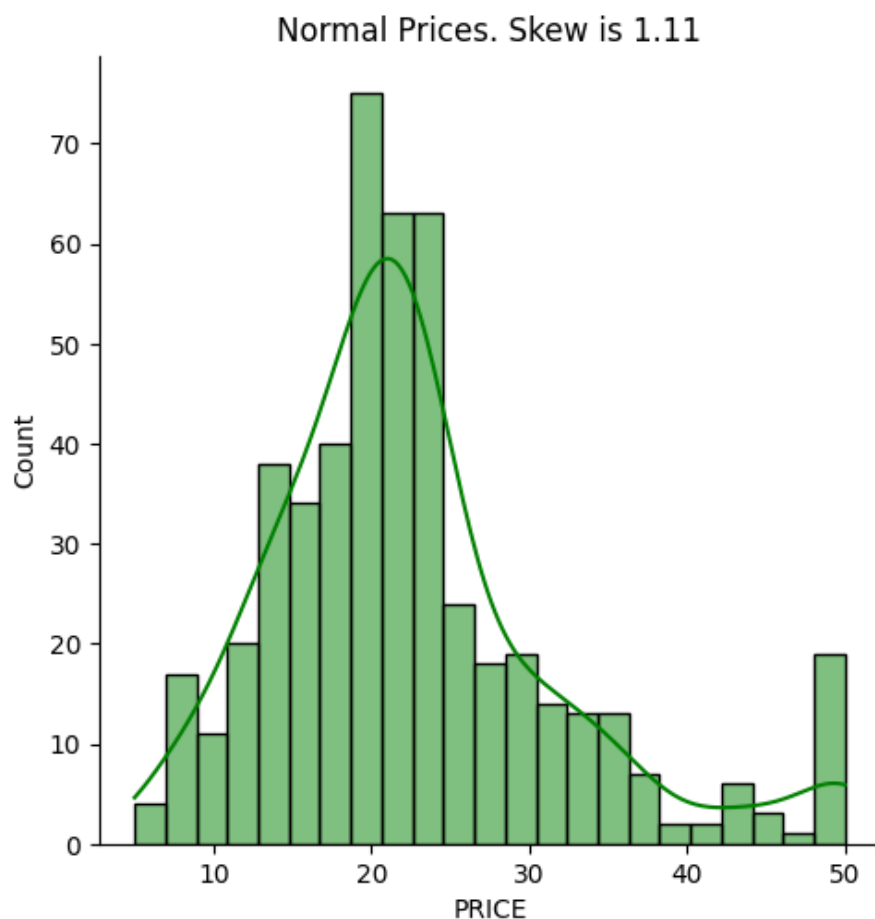**Let's try a data transformation approach.**

**Challenge**

Investigate if the target `data['PRICE']` could be a suitable candidate for a log transformation.

- **Use Seaborn's `.displot()` to show a histogram and KDE of the price data.**
- **Calculate the skew of that distribution.**
- **Use NumPy's `log()` function to create a Series that has the log prices**
- **Plot the log prices using Seaborn's `.displot()` and calculate the skew.**
- **Which distribution has a skew that's closer to zero?**

In [66]:

```
tgt_skew = data['PRICE'].skew()
sns.displot(data['PRICE'], kde='kde', color='green')
plt.title(f'Normal Prices. Skew is {tgt_skew:.3}')
plt.show()
```
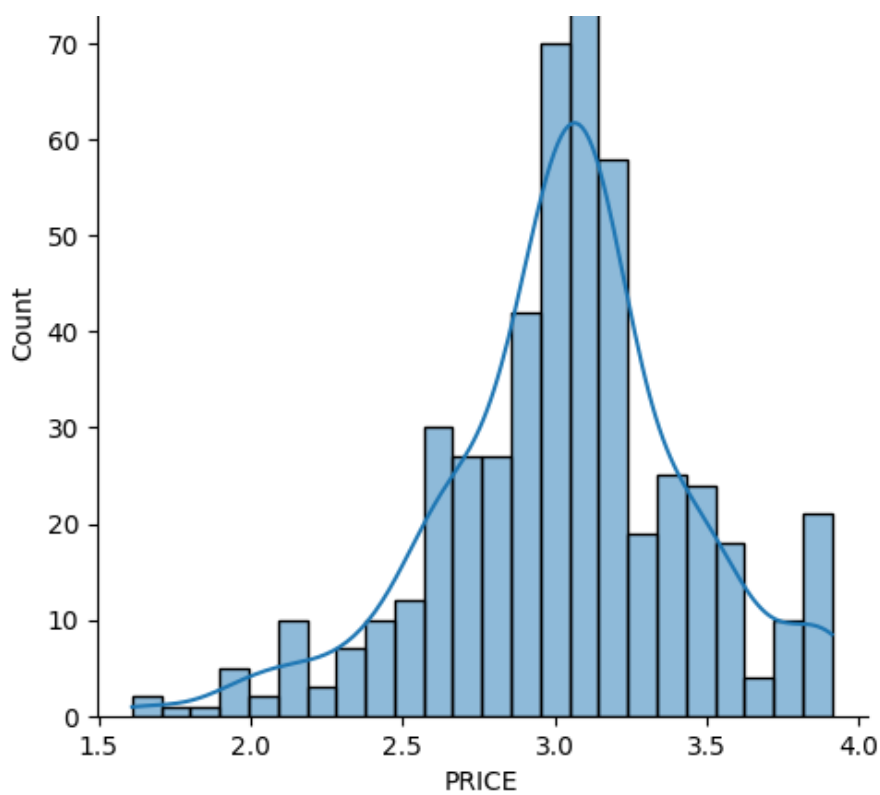


In [67]:

```
y_log = np.log(data['PRICE'])
sns.displot(y_log, kde=True)
plt.title(f'Log Prices. Skew is {y_log.skew():.3}')
plt.show()
```
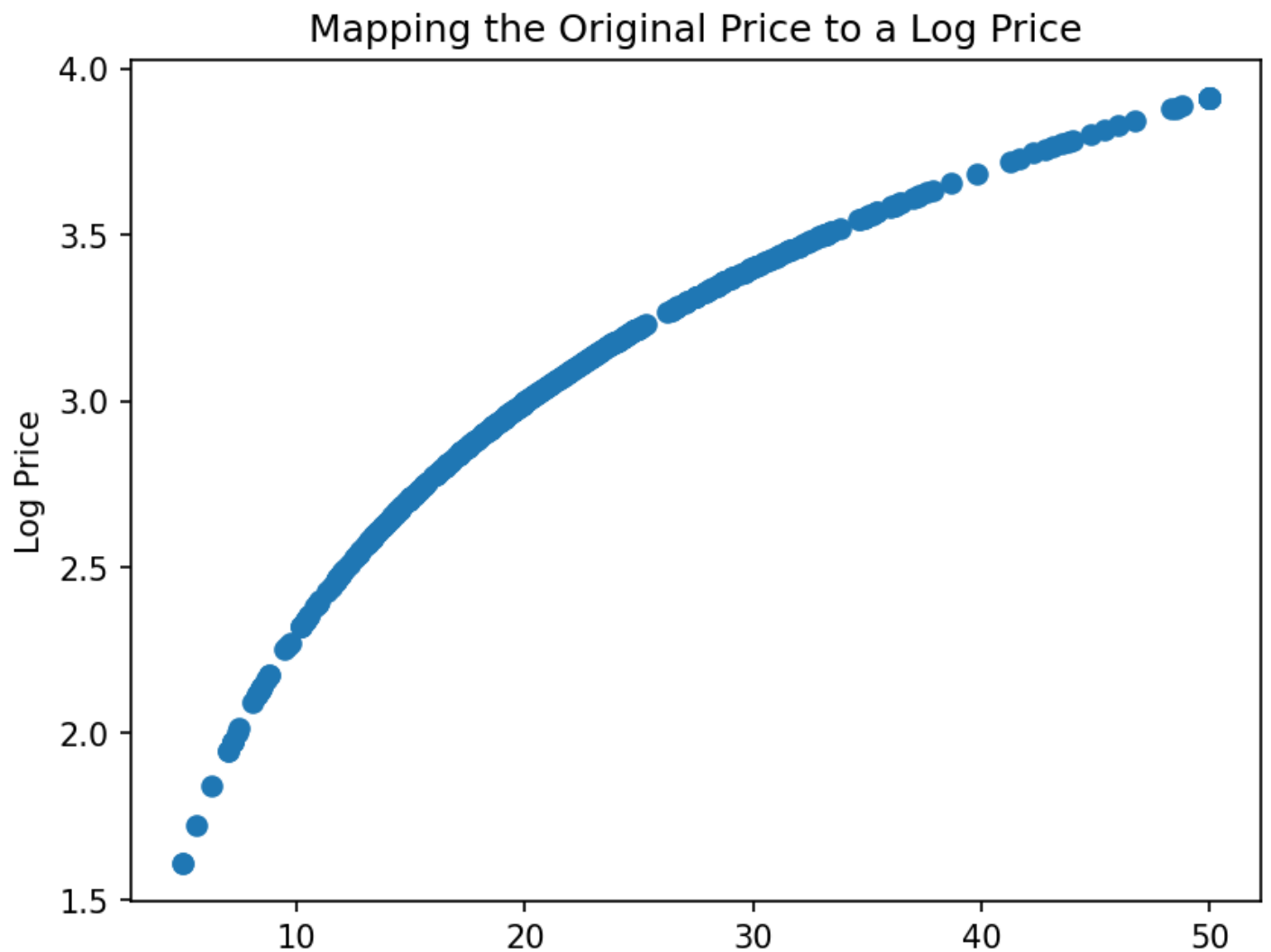
```
plt.figure(dpi=150)
plt.scatter(data.PRICE, np.log(data.PRICE))

plt.title('Mapping the Original Price to a Log Price')
plt.ylabel('Log Price')
plt.xlabel('Actual $ Price in 000s')
plt.show()
```

**How does the log transformation work?**

Using a log transformation does not affect every price equally. Large prices are affected more than smaller prices in the dataset. Here's how the prices are "compressed" by the log transformation:
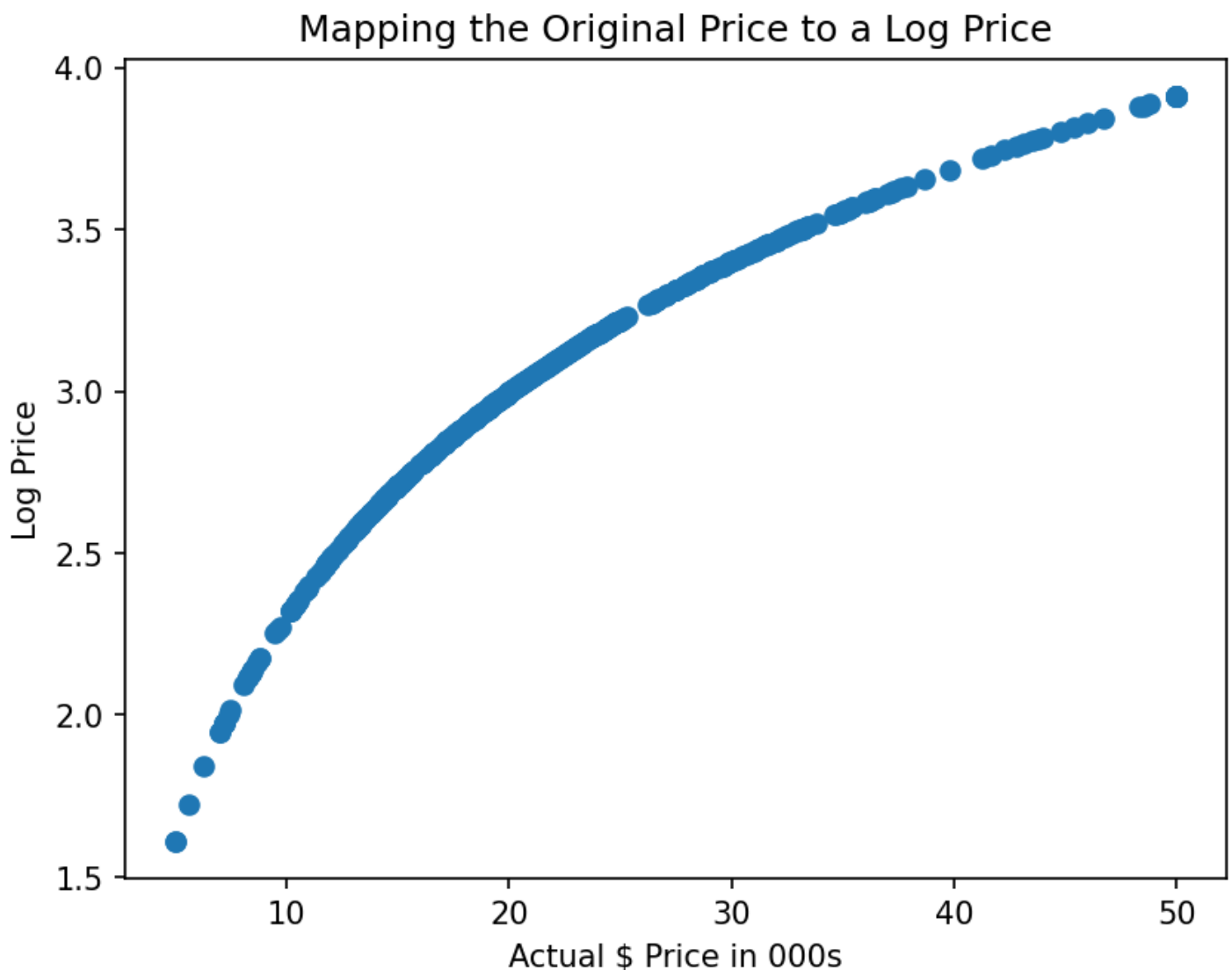
| If PRICE = 7 ... | -5.05 → | ln(7) = 1.95 |
|---|---|---|
| If PRICE = 50 ... | -46.09 → | ln(50) = 3.91 |

We can see this when we plot the actual prices against the (transformed) log prices.

In [55]:

```python
plt.figure(dpi=150)
plt.scatter(data.PRICE, np.log(data.PRICE))

plt.title('Mapping the Original Price to a Log Price')
plt.ylabel('Log Price')
plt.xlabel('Actual $ Price in 000s')
plt.show()
```



Mapping the Original Price to a Log Price

# Regression using Log Prices

Using log prices instead, our model has changed to:

$$\log(\hat{PRICE}) = \theta_0$$
$$+ \theta_1 RM + \theta_2 NOX$$
$$+ \theta_3 DIS + \theta_4 CHAS +.$$
$$..+\theta_{13} LSTAT$$

**Challenge:**

- Use `train_test_split()` with the same random state as before to make the results comparable.
- Run a second regression, but this time use the transformed target data.
- What is the r-squared of the regression on the training data?
- Have we improved the fit of our model compared to before based on this measure?

In [70]:

```python
new_target = np.log(data['PRICE']) # Use log prices
features = data.drop('PRICE', axis=1)

X_train, X_test, log_y_train, log_y_test = train_test_split(features,
                                                             new_target,
                                                             test_size=0.2,
                                                             random_state=10)

log_regr = LinearRegression()
log_regr.fit(X_train, log_y_train)
log_rsquared = log_regr.score(X_train, log_y_train)

log_predictions = log_regr.predict(X_train)
log_residuals = (log_y_train - log_predictions)

print(f'Training data r-squared: {log_rsquared:.2}')
```

```
Training data r-squared: 0.79
```

In [55]:

# Evaluating Coefficients with Log Prices

**Challenge:** Print out the coefficients of the new regression model.

- Do the coefficients still have the expected sign?
- Is being next to the river a positive based on the data?
- How does the quality of the schools affect property prices? What happens to prices as there are more students per teacher?

**Hint: Use a DataFrame to make the output look pretty.**

In [71]:

```python
df_coef = pd.DataFrame(data=log_regr.coef_, index=X_train.columns, columns=['coef'])
df_coef
```

Out[71]:

|       | coef  |
|-------|-------|
| CRIM  | -0.01 |
| ZN    | 0.00  |
| INDUS | 0.00  |
| CHAS  | 0.08  |

| | |
|---|---|
| NOX | 0.79 |
| RM | 0.07 |
| AGE | 0.00 |
| DIS | -0.05 |
| RAD | 0.01 |
| TAX | -0.00 |
| PTRATIO | -0.03 |
| B | 0.00 |
| LSTAT | -0.03 |

In [55]:

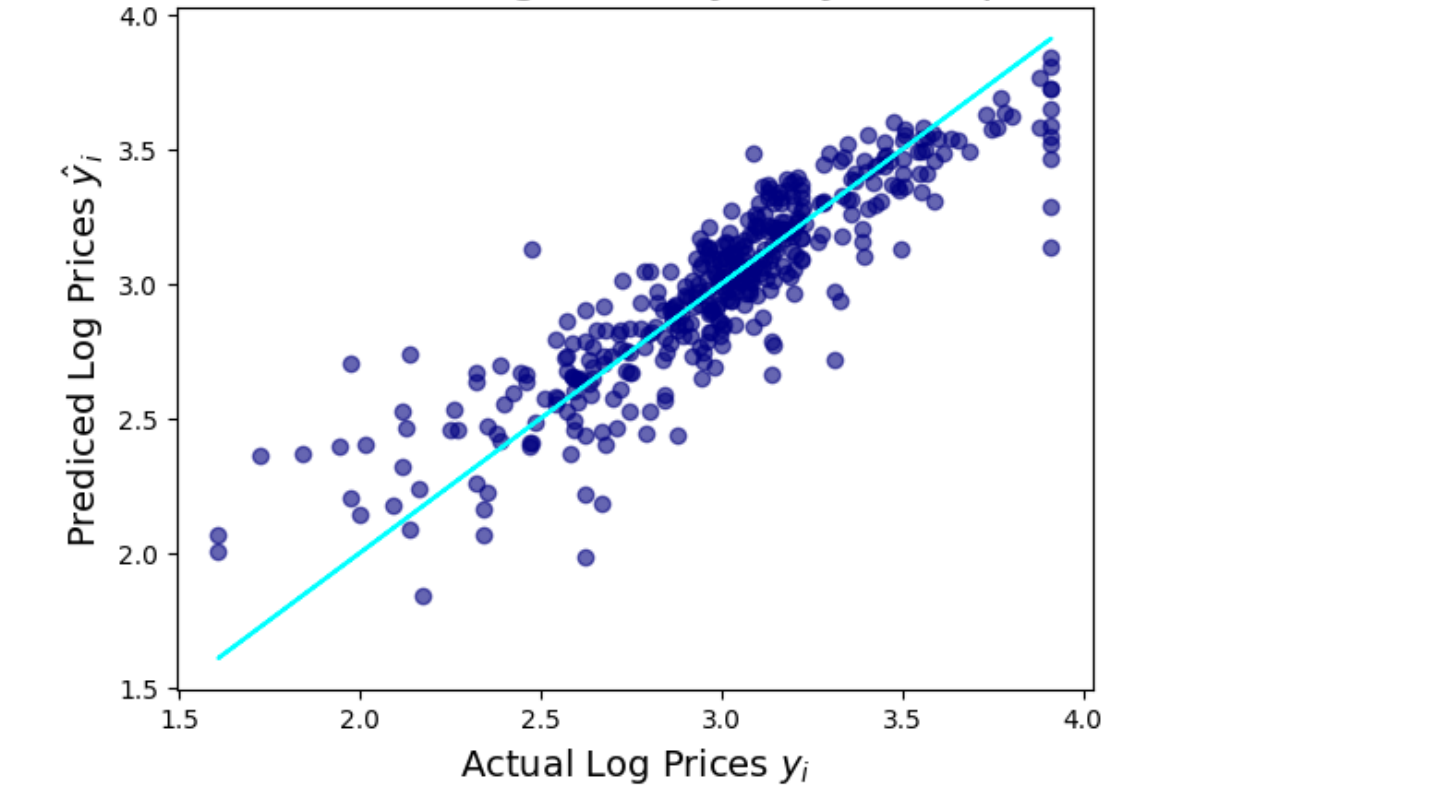## Regression with Log Prices & Residual Plots

**Challenge:**

- Copy-paste the cell where you've created scatter plots of the actual versus the predicted home prices as well as the residuals versus the predicted values.
- Add 2 more plots to the cell so that you can compare the regression outcomes with the log prices side by side.
- Use `indigo` as the colour for the original regression and `navy` for the color using log prices.

In [72]:

```python
# Graph of Actual vs. Predicted Log Prices
plt.scatter(x=log_y_train, y=log_predictions, c='navy', alpha=0.6)
plt.plot(log_y_train, log_y_train, color='cyan')
plt.title(f'Actual vs Predicted Log Prices: $y _i$ vs $\hat y_i$ (R-Squared {log_rsquared
:.2})', fontsize=17)
plt.xlabel('Actual Log Prices $y _i$', fontsize=14)
plt.ylabel('Prediced Log Prices $\hat y _i$', fontsize=14)
plt.show()
```
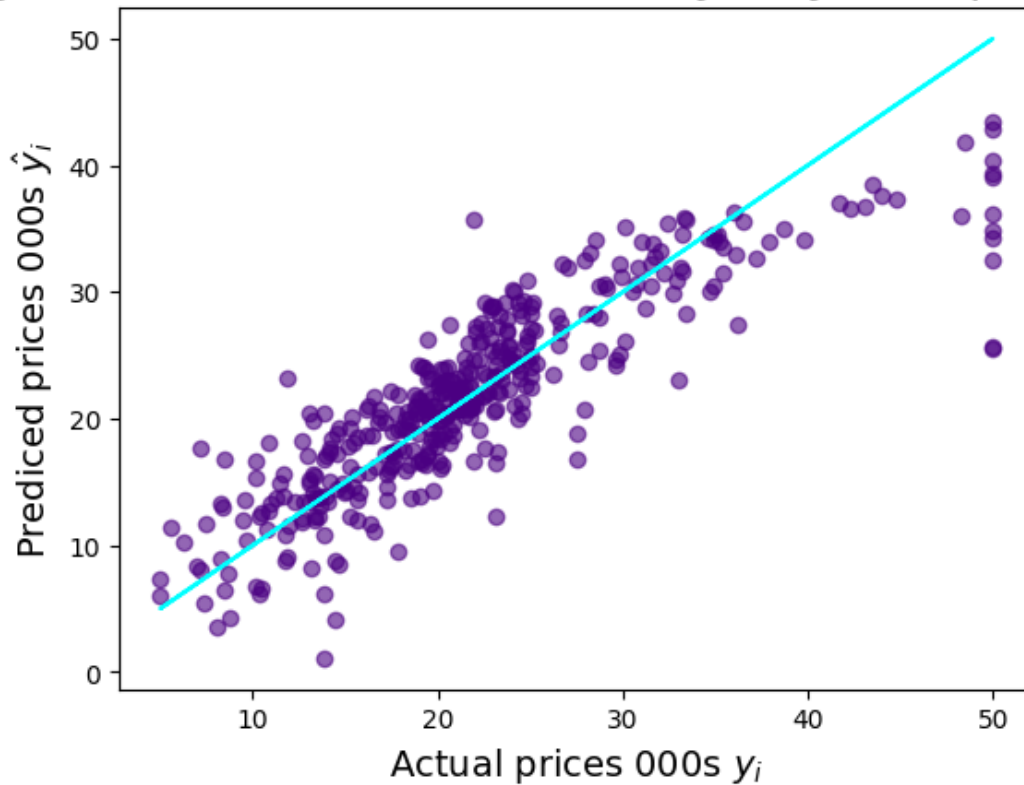


Actual vs Predicted Log Prices: $y_i$ vs $\hat{y}_i$ (R-Squared 0.79)

In [73]:

```python
# Original Regression of Actual vs. Predicted Prices
plt.scatter(x=y_train, y=predicted_vals, c='indigo', alpha=0.6)
plt.plot(y_train, y_train, color='cyan')
plt.title(f'Original Actual vs Predicted Prices: $y _i$ vs $\hat y_i$ (R-Squared {rsquare
d:.3})', fontsize=17)
plt.xlabel('Actual prices 000s $y _i$', fontsize=14)
plt.ylabel('Prediced prices 000s $\hat y _i$', fontsize=14)
plt.show()
```
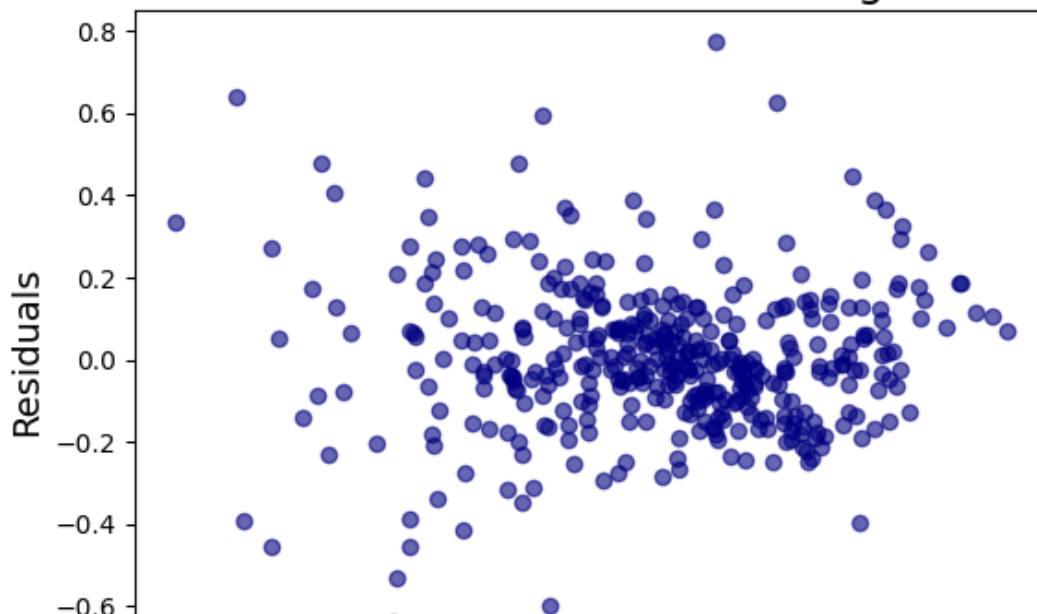
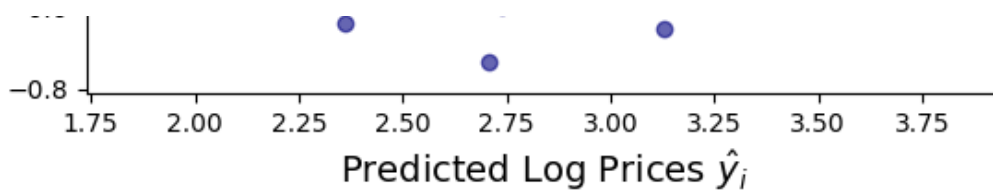Original Actual vs Predicted Prices: $y_i$ vs $\hat y_i$ (R-Squared 0.75)



In [74]:

```python
# Residuals vs Predicted values (Log prices)
plt.scatter(x=log_predictions, y=log_residuals, c='navy', alpha=0.6)
plt.title('Residuals vs Fitted Values for Log Prices', fontsize=17)
plt.xlabel('Predicted Log Prices $\hat y _i$', fontsize=14)
plt.ylabel('Residuals', fontsize=14)
plt.show()
```
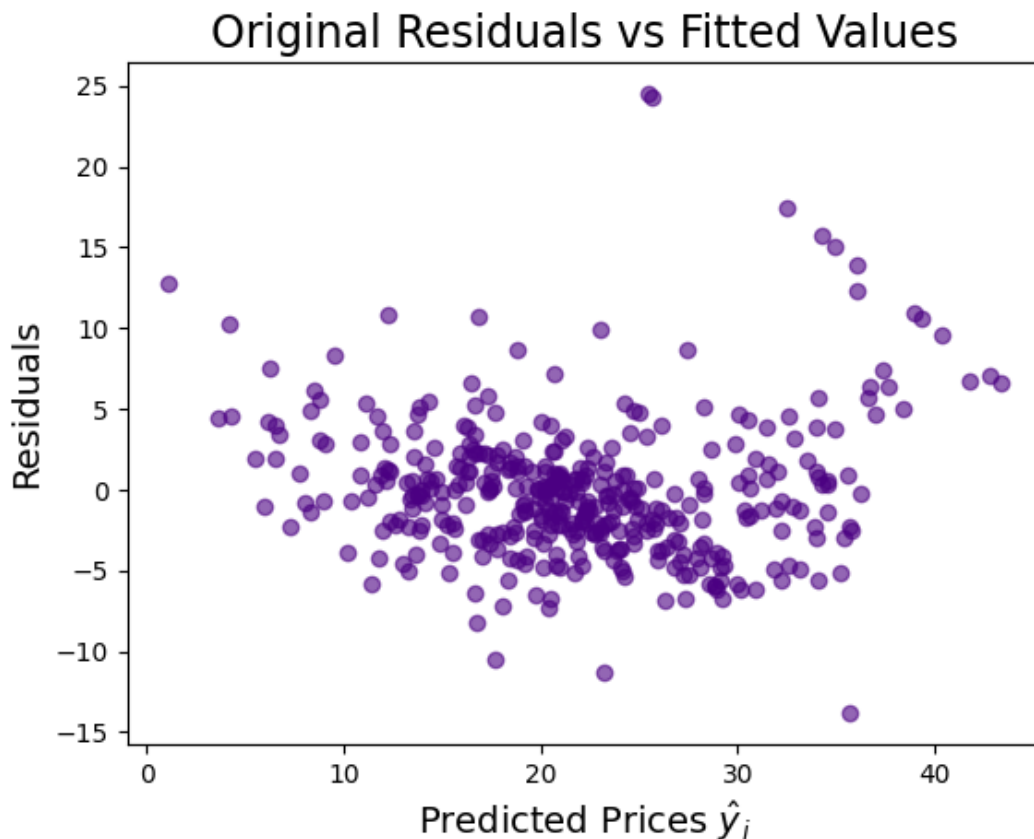
Residuals vs Fitted Values for Log Prices

−0.8

1.75   2.00   2.25   2.50   2.75   3.00   3.25   3.50   3.75

Predicted Log Prices $\hat y_i$

In [75]:

```python
# Residuals vs Predicted values
plt.scatter(x=predicted_vals, y=residuals, c='indigo', alpha=0.6)
plt.title('Original Residuals vs Fitted Values', fontsize=17)
plt.xlabel('Predicted Prices $\hat y _i$', fontsize=14)
plt.ylabel('Residuals', fontsize=14)
plt.show()
```



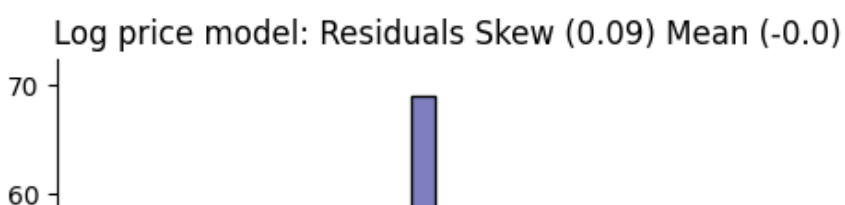## Original Residuals vs Fitted Values

**Challenge:**

**Calculate the mean and the skew for the residuals using log prices. Are the mean and skew closer to 0 for the regression using log prices?**
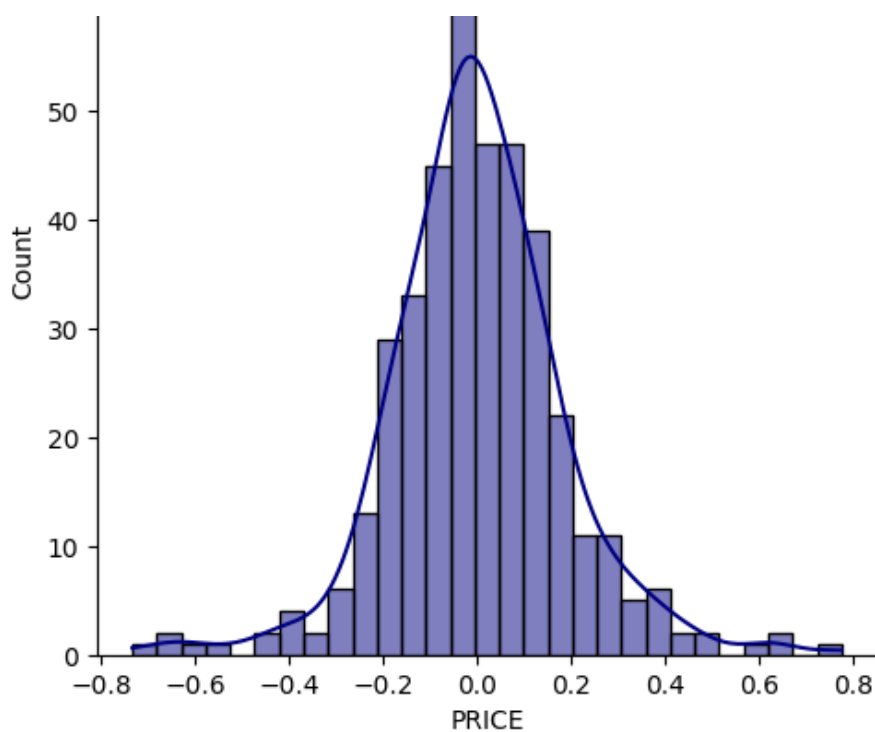
In [76]:

```python
# Distribution of Residuals (log prices) - checking for normality
log_resid_mean = round(log_residuals.mean(), 2)
log_resid_skew = round(log_residuals.skew(), 2)

sns.displot(log_residuals, kde=True, color='navy')
plt.title(f'Log price model: Residuals Skew ({log_resid_skew}) Mean ({log_resid_mean})')
plt.show()

sns.displot(residuals, kde=True, color='indigo')
plt.title(f'Original model: Residuals Skew ({resid_skew}) Mean ({resid_mean})')
plt.show()
```
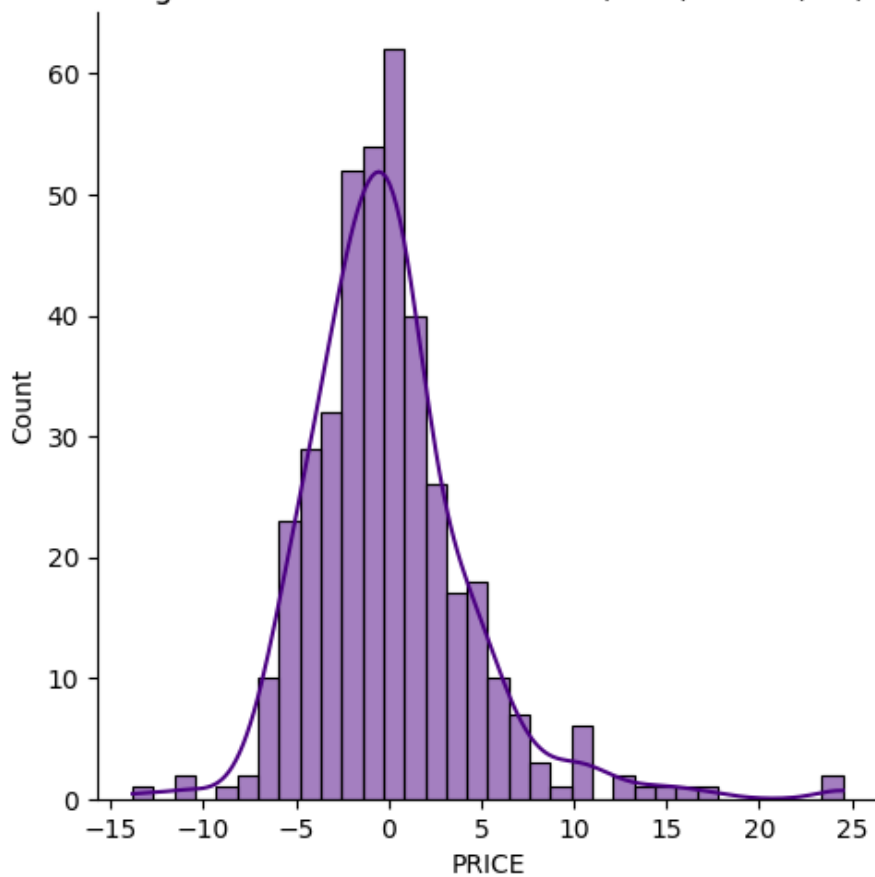
Log price model: Residuals Skew (0.09) Mean (-0.0)

70

60

Original model: Residuals Skew (1.46) Mean (0.0)



In [55]:

## Compare Out of Sample Performance

The *real* test is how our model performs on data that it has not "seen" yet. This is where our `X_test` comes in.

**Challenge**

Compare the r-squared of the two models on the test dataset. Which model does better? Is the r-squared higher or lower than for the training dataset? Why?

```
print(f'Original Model Test Data r-squared: {regr.score(X_test, y_test):.2}')
print(f'Log Model Test Data r-squared: {log_regr.score(X_test, log_y_test):.2}')
```

```
Original Model Test Data r-squared: 0.67
Log Model Test Data r-squared: 0.74
```

In [55]:

# Predict a Property's Value using the Regression Coefficients

**Our preferred model now has an equation that looks like this:**

$$\log(\hat{PRICE}) = \theta_0$$
$$+ \theta_1 RM + \theta_2 NOX$$
$$+ \theta_3 DIS + \theta_4 CHAS + ..$$
$$. + \theta_{13} LSTAT$$

**The average property has the mean value for all its charactistics:**

In [78]:

```
# Starting Point: Average Values in the Dataset
features = data.drop(['PRICE'], axis=1)
average_vals = features.mean().values
property_stats = pd.DataFrame(data=average_vals.reshape(1, len(features.columns)),
                              columns=features.columns)
property_stats
```

Out[78]:

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3.61 | 11.36 | 11.14 | 0.07 | 0.55 | 6.28 | 68.57 | 3.80 | 9.55 | 408.24 | 18.46 | 356.67 | 12.65 |

**Challenge**

**Predict how much the average property is worth using the stats above. What is the log price estimate and what is the dollar estimate? You'll have to [reverse the log transformation with](# ) `.exp()` to find the dollar value.**

In [79]:

```
# Make prediction
log_estimate = log_regr.predict(property_stats)[0]
print(f'The log price estimate is ${log_estimate:.3}')

# Convert Log Prices to Acutal Dollar Values
dollar_est = np.e**log_estimate * 1000
# or use
dollar_est = np.exp(log_estimate) * 1000
print(f'The property is estimated to be worth ${dollar_est:.6}')
```

```
The log price estimate is $3.03
The property is estimated to be worth $20703.2
```

In [56]:

**Challenge**

**Keeping the average values for CRIM, RAD, INDUS and others, value a property with the following characteristics:**

In [57]:

```
# Define Property Characteristics
next_to_river = True
nr_rooms = 8
students_per_classroom = 20
distance_to_town = 5
pollution = data.NOX.quantile(q=0.75) # high
amount_of_poverty =  data.LSTAT.quantile(q=0.25) # low
```

In [58]:

```
# Solution:
```

In [58]: