

```
In [40]: # Initialize Otter
import otter
grader = otter.Notebook("lab6-regression.ipynb")
```

```
In [41]: import numpy as np
import math
import pandas as pd
import altair as alt
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
from sklearn.preprocessing import add_dummy_feature

alt.data_transformers.disable_max_rows()
```

```
Out[41]: DataTransformerRegistry.enable('default')
```

# Lab 6: Regression

This lab covers the nuts and bolts of fitting linear models using the `sklearn.linear_model` module. The linear model expresses a response variable,  $y$ , as a linear function of  $p - 1$  explanatory variables  $x_1, \dots, x_{p-1}$  and a random error  $\epsilon$ . Its general form is:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1} + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

Usually, the response and explanatory variables and error term are indexed by observation  $i = 1, \dots, n$  so that the model describes a dataset comprising  $n$  values of each variable:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{i,p-1} + \epsilon_i \quad \begin{cases} \epsilon_i \sim N(0, \sigma^2) \\ i = 1, \dots, n \end{cases}$$

Because the indices get confusing to keep track of, it is much easier to express the model in matrix form as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1,p-1} \\ 1 & x_{21} & \cdots & x_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{n,p-1} \end{bmatrix}_{n \times p} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix}_{p \times 1} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

**Fitting** a model of this form means *estimating the parameters*  $\beta_0, \beta_1, \dots, \beta_{p-1}$  and  $\sigma^2$  from a set of data.

- The ordinary least squares estimates of  $\boldsymbol{\beta}$ , which are best under most circumstances, are

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- The error variance  $\sigma^2$  can be estimated by

$$\hat{\sigma}^2 = \frac{1}{n - p - 1} \left( \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \right)' \left( \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \right)$$

When fitting a linear model, it is also of interest to quantify uncertainty by estimating the variability of  $\hat{\boldsymbol{\beta}}$  and measure overall quality of fit. This lab illustrates that process and the computations involved.

## Objectives

In this lab, you'll learn how to use the `sklearn.linear_model` module to:

- compute OLS estimates;
- calculate fitted values and residuals;
- compute the error variance estimate.

In addition, you'll see how to calculate:

- the variance-covariance matrix of  $\hat{\boldsymbol{\beta}}$ , which quantifies the variability of model estimates;
- standard errors for each model estimate;
- the proportion of variation captured by a linear model.

Throughout you'll use simple visualizations to help make the connection between fitted models and the aspects of a dataset that model features describe. The lab activity proceeds through the following sections:

- 0. Exploratory plots for regression analysis
- 1. Simple linear regression (single explanatory variable)
  - Model fitting

- Uncertainty quantification
  - Model visualization
2. Multiple linear regression (several explanatory variables)
- Data preprocessing: encoding categorical variables
  - Interaction terms
  - Model fitting
  - Uncertainty quantification
  - Model visualization

## Data: fertility rates and meaasures of development

By way of data, you'll work with country indicators, total fertility rates, and gender indicators for a selection of countries in 2018, and explore the decline in fertility rates associated with developed nations. This has been in the news lately due to preliminary data from the U.S. 2020 census indicating significant [population growth decline in the United States](#). If the topic interests you, you can read more about perspectives and existing data in this [Our World in Data](#) article.

The data are stored in separate .csv files imported below:

```
In [42]: fertility = pd.read_csv('data/fertility.csv')
country = pd.read_csv('data/country-indicators.csv')
gender = pd.read_csv('data/gender-data.csv')
```

The variables you'll work with in this portion are the following:

Dataset	Name	Variable	Units
fertility	fertility_total	National fertility rate	Average number of children per woman
country	hdi	Human development index	Index between 0 and 1 (0 is lowest, 1 is highest)
gender	edu_expected_yrs_f	Expected years of education for adult women	Years

Because the variables of interest are stored in three separate dataframes, you'll first need to extract them and merge by country.

```
In [43]: # slice variables of interest
fertility_sub = fertility.loc[:, ['Country', 'fertility_total']]
gender_sub = gender.loc[:, ['educ_expected_yrs_f', 'Country']]
country_sub = country.loc[:, ['Country', 'hdi']]

# merge variables of interest
reg_data = pd.merge(
    fertility_sub,
    gender_sub,
    on = 'Country',
    how = 'inner'
).merge(
    country_sub,
    on = 'Country',
    how = 'left'
).set_index('Country').dropna()

# preview
reg_data.head(4)
```

```
Out[43]:
```

	fertility_total	educ_expected_yrs_f	hdi
Country			
Afghanistan	4.473	6.795722	0.509
Albania	1.617	13.201755	0.792
Algeria	3.023	12.108990	0.746
Angola	5.519	6.973901	0.582

## 0. Exploratory plots

A preliminary step in regression analysis is typically data exploration through scatterplots. This relies on skills you've already developed -- you may find it helpful to refer to Lab 3 (visualization) throughout this section.

### Q0 (a) Simple scatterplot

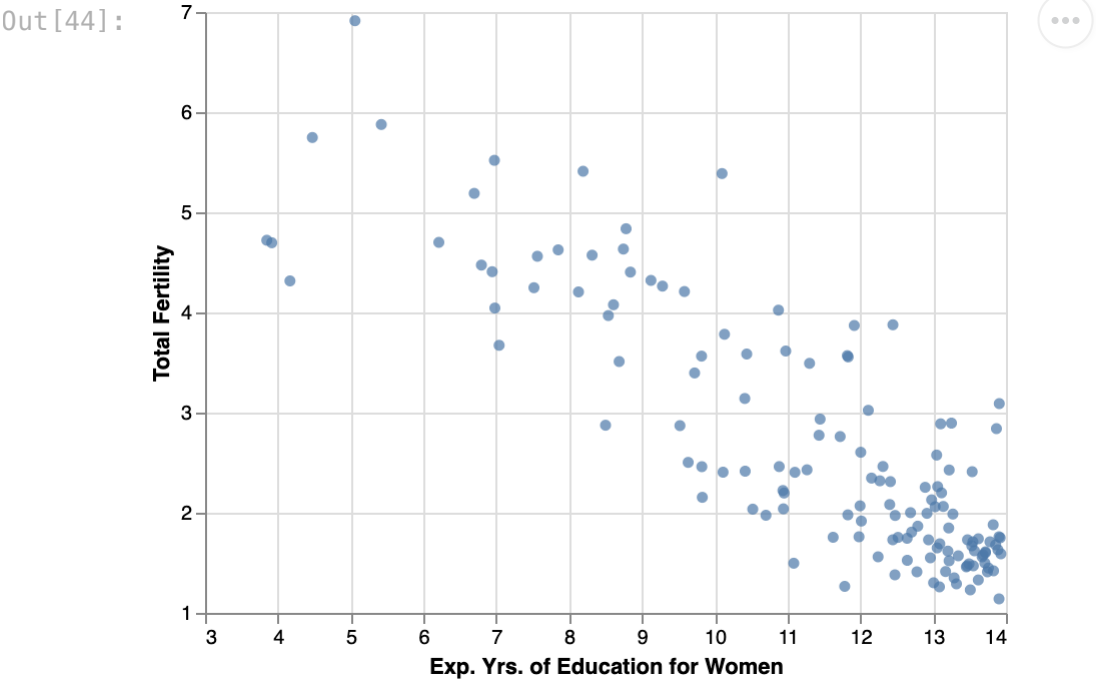
Construct a scatterplot of total fertility against expected years of education for women. Label the axes 'Fertility rate' and 'Expected years of education for women'. Store this plot as `simple_scatter` and display the graphic.

(Remark: be sure to include `scale = alt.Scale(zero = False)` in the axis specification so that your plot does not have extra whitespace.)

```
In [44]: #solution
simple_scatter = alt.Chart(reg_data).mark_circle().encode(
```

```
x = alt.X('educ_expected_yrs_f', title = 'Exp. Yrs. of Education for Women', scale = alt.Scale(zero = False)),
y = alt.Y('fertility_total', title = 'Total Fertility', scale = alt.Scale(zero = False))
)

simple_scatter
```



The cell below coerces the human development index to a high-med-low factor (categorical variable) using `pd.qcut(...)` . The cutoffs for the factor levels are determined by sample quantiles.

```
In [45]: # labels for the factor
hdi_labels = ['1_low', '2_medium', '3_high']

# cut hdi based on quantiles and store as new variable
reg_data['hdi_fac'] = pd.qcut(reg_data.hdi,
                             q = 3,
                             labels = hdi_labels)

# preview
reg_data.head(3)
```

Out[45]:

	fertility_total	educ_expected_yrs_f	hdi	hdi_fac
Country				
Afghanistan	4.473	6.795722	0.509	1_low
Albania	1.617	13.201755	0.792	2_medium
Algeria	3.023	12.108990	0.746	2_medium

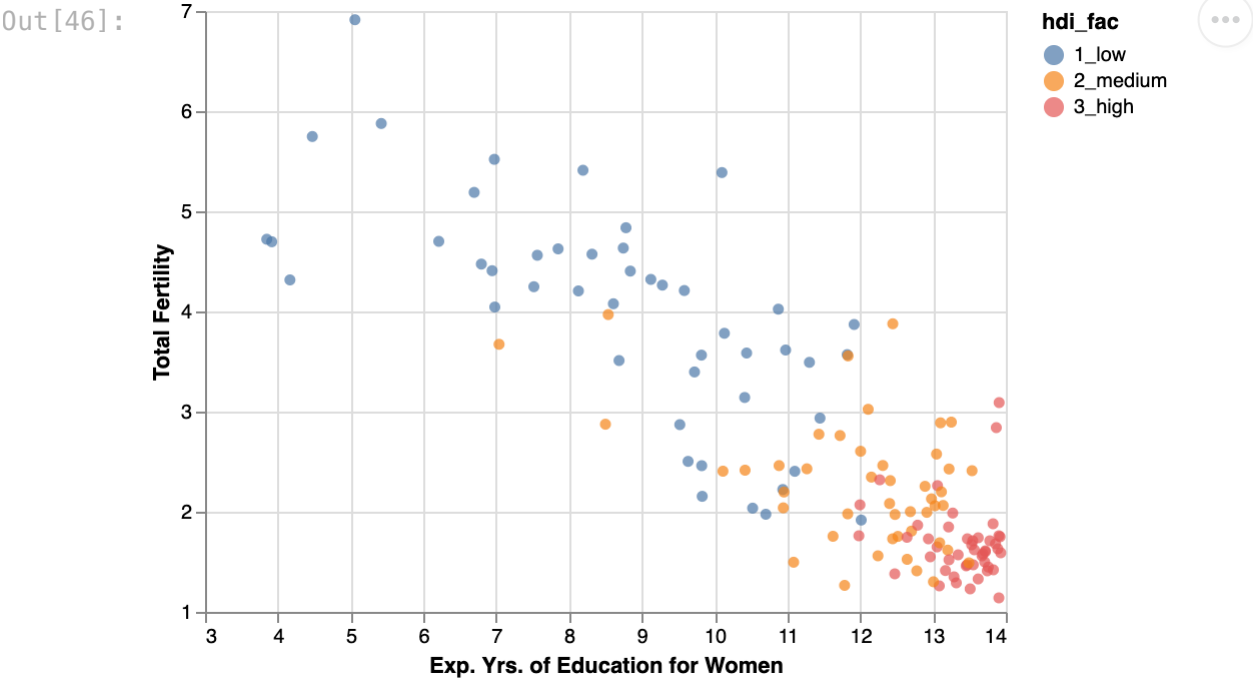
## Q0 (b) Scatter by HDI level

### (i) Add HDI level

Modify your plot from Q0 (a) so that points are colored according to the corresponding country's level of human development. Make sure you choose appropriate labels for your axes and plot. Store this plot as `scatter` and display the graphic.

```
In [46]: scatter = simple_scatter.encode(color = 'hdi_fac:N')

scatter
```



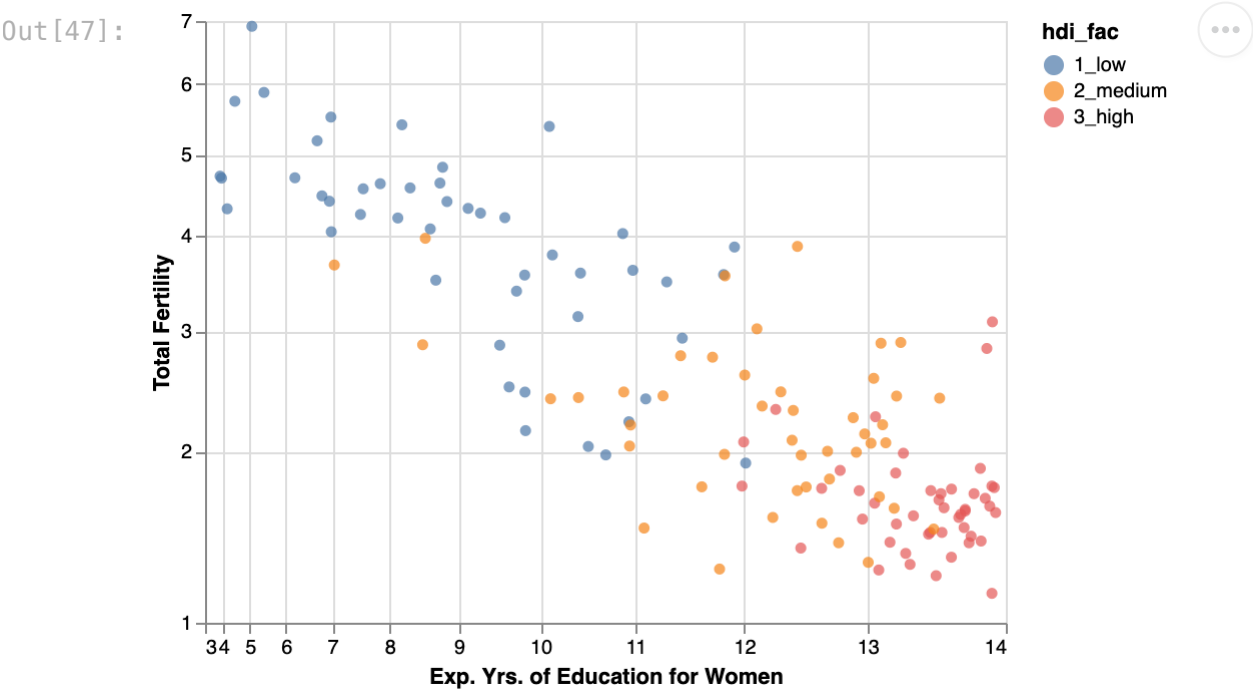
### (ii) Adjust axis scales

Notice that the scatter is clumped together densely in the lower right corner. To better see what's happening in that region, modify your codes for constructing `scatter` so that your axis specifications for both axes include `alt.Scale(zero = False, type = 'pow', exponent = ...)` .

Make sure you choose appropriate labels for your axes and plot.

This will display the axes on a power-transformed scale with an exponent of your choosing. Use an exponent around 2.5 for the education axis, and an exponent around 1/3 for the fertility axis. Store the plot as `scatter_scaled` and display the result.

```
In [47]: scatter_scaled = scatter.encode(  
    x = alt.X('educ_expected_yrs_f', title = 'Exp. Yrs. of Education for Women',  
        scale = alt.Scale(zero = False, type = 'pow', exponent = 2.5)),  
    y = alt.Y('fertility_total', title = 'Total Fertility',  
        scale = alt.Scale(zero = False, type = 'pow', exponent = .333))  
)  
  
scatter_scaled
```

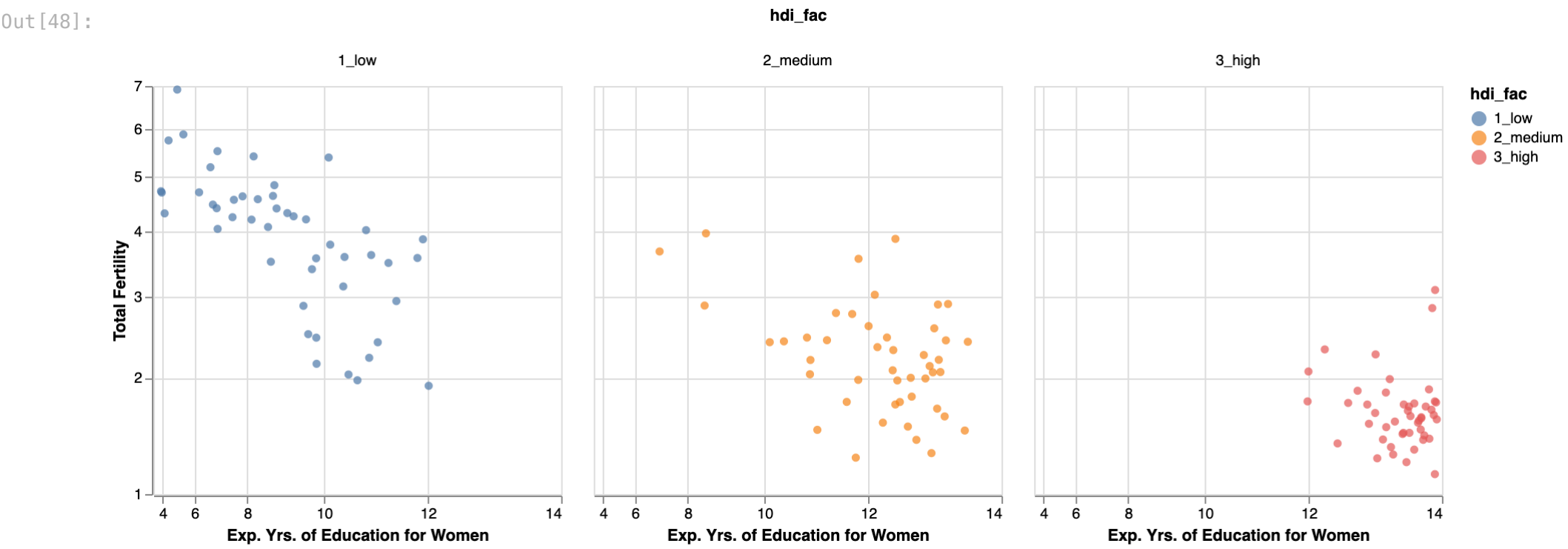


Q0 (c) Interpret your graphic

Does the negative association (downward trend in the scatter) seem to hold steady among countries at each level of human development? Try to imagine drawing a line to fit the points of each color, and ask yourself whether the slopes would be the same.

i. Facet by human development ( `hdi_fac` level) and store the result in `faceted_hdi` so that you can see the pattern for each development level apart from the other levels.

```
In [48]: # facet by hdi level  
faceted_hdi = scatter_scaled.properties(  
    width = 275,  
    height = 275  
)  
.facet(column = 'hdi_fac')  
  
faceted_hdi
```



Does the negative association (downward trend in the scatter) seem to hold steady among countries at each level of human development?

Yes the negative association seems to hold steady among the countries at each level of human development.

Examining individual points

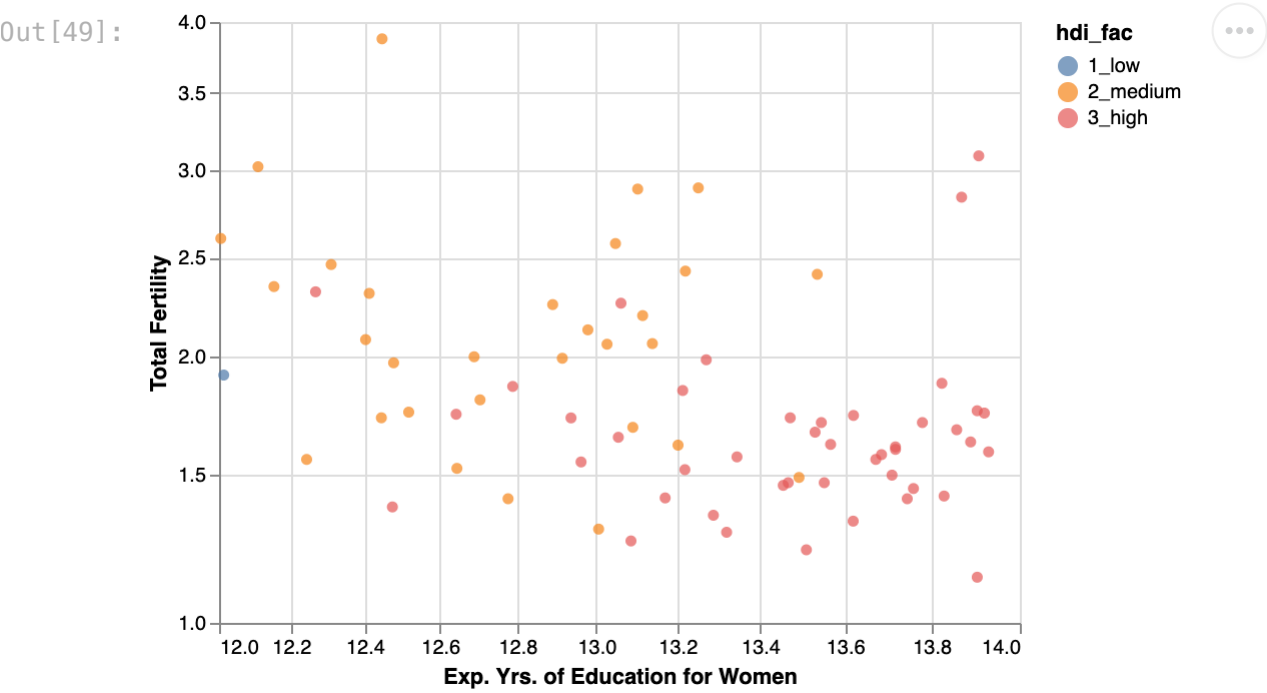
Often it can be helpful to graphically examine individual points. Especially so if you notice outliers. In this dataset there aren't any apparent outliers, but there are a few points that seem a little farther from the scatter than others, so we can practice on those.

Let's zoom in on the lower-right region of the scatterplot; the cell below applies an Altair filter transform to show only data for which expected education exceeds 12.

Please be sure you've created `scatter_scaled` in Q0 (b) (ii) before running this cell.

```
In [49]: # zoom in
scatter_highed = scatter_scaled.transform_filter(
    alt.FieldGTPredicate(field = 'educ_expected_yrs_f', gt = 12)
)

scatter_highed
```



Notice that there are a few points that stand apart from the main scatter. Which countries do those points represent? It's not too tricky to pick them out by filtering and taking maxima or minima. For example, the country with the lowest fertility rate is:

```
In [50]: reg_data[reg_data.educ_expected_yrs_f > 12].fertility_total.idxmin()

Out[50]: 'Singapore'
```

Q0 (e) Labeling specific countries

Here you'll look at a few more specific countries, and then label them on the plot. Consider countries whose expected education exceeds 12.

(i) Find the country among the high-HDI nations with the highest total fertility and store it as `high_high_fert`.

```
In [51]: high_high_fert = reg_data[(reg_data.educ_expected_yrs_f > 12) & (reg_data.hdi_fac == '3_high')].fertility_total.idxmax()

print(high_high_fert)

Israel
```

(ii) Find the country among the medium-HDI nations with the highest total fertility and store it as `med_high_fert`.

```
In [52]: med_high_fert = reg_data[(reg_data.educ_expected_yrs_f > 12) & (reg_data.hdi_fac == '2_medium')].fertility_total.idxmax()

print(med_high_fert)

Samoa
```

(iii) Identify one other country that calls your attention.

```
In [53]: #BEGIN SOLUTION NO PROMPT
#potential solution
low_high_fert = reg_data[(reg_data.hdi_fac == '1_low') & (reg_data.educ_expected_yrs_f > 12)].fertility_total.idxmax()

print(low_high_fert)
#END SOLUTION

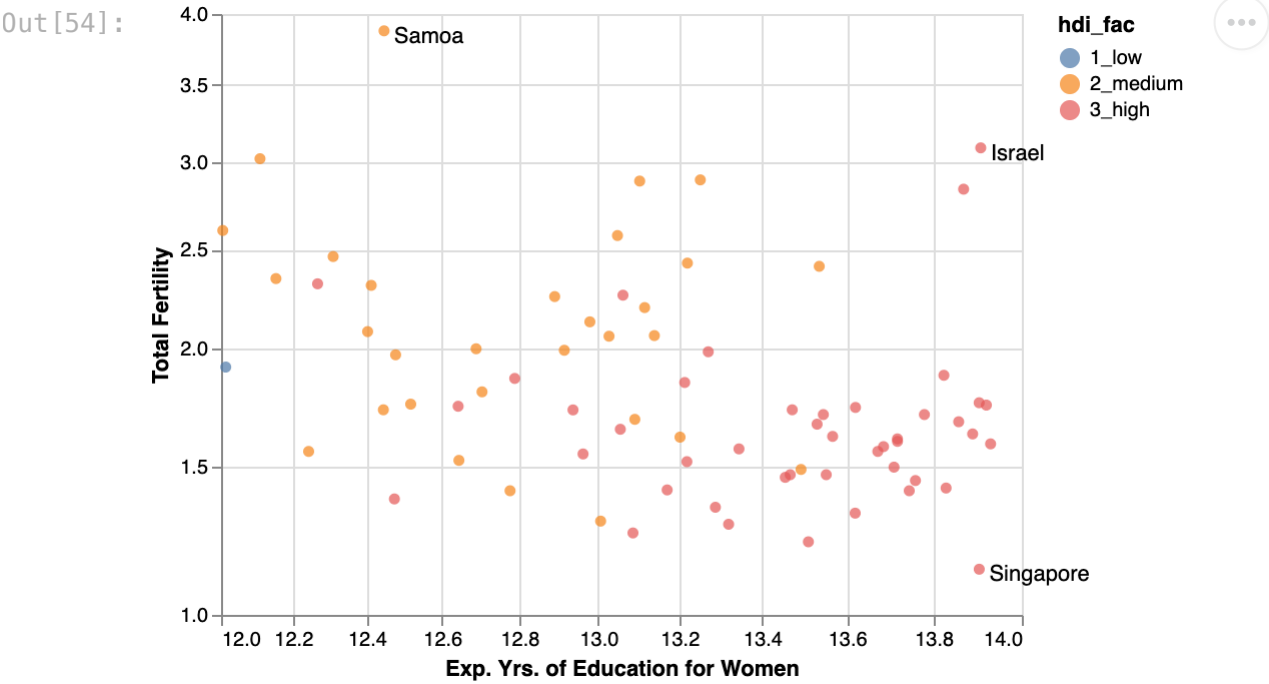
Nepal
```

(iv) Modify the cell below to label the three countries identified above.

```
In [54]: # slice rows for identified countries
point_labels = reg_data.loc[['Singapore', 'Israel', 'Samoa']].reset_index() # edit here

# plot text labels
labels = alt.Chart(point_labels).mark_text(align = 'left', dx = 5, dy = 3).encode(
    x = 'educ_expected_yrs_f',
    y = 'fertility_total',
    text = 'Country'
)

# layer over scatter
scatter_highed + labels
```



## 1. Simple linear regression

In this part you'll fit a simple linear model regressing fertility on education.

First we'll need to store the quantities -- the response and explanatory variables -- needed for model fitting in the proper format. Recall that the linear model in matrix form is:

$$\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}}_{\boldsymbol{\epsilon}}$$

Notice that the explanatory variable matrix  $\mathbf{X}$  includes a column of ones for the intercept. So the quantities needed are:

- $\mathbf{y}$ , a one-dimensional array of the total fertility rates for each country; and
- $\mathbf{X}$ , a two-dimensional array with a column of ones (intercept) and a column of the expected years of education for women (explanatory variable).

The cell below prepares the explanatory variable matrix  $\mathbf{X}$  and response vector  $\mathbf{y}$  in the following steps:

1. Store the total fertility rates as an array named `y`.
2. Slice the `educ_expected_yrs_f` variable from the `reg_data` dataframe, resulting in a dataframe with one column stored as `x_slr_df`.
3. Append a column of ones to `x_slr_df` using `add_dummy_feature`, resulting in a 139-by-2 array stored as `x_slr`.

```
In [55]: # store response variable as array
y = reg_data.fertility_total

# store explanatory variable
x_slr_df = reg_data.loc[:, ['educ_expected_yrs_f']]

# add intercept column
x_slr = add_dummy_feature(x_slr_df, value = 1)

# print first five rows
x_slr[0:5, :]
```

```
Out[55]: array([[ 1.          ,  6.79572248],
 [ 1.          , 13.20175457],
 [ 1.          , 12.10898972],
 [ 1.          ,  6.97390127],
 [ 1.          , 12.91444111]])
```

### Model fitting

'Fitting' a model refers to computing estimates; it is typical practice to report estimated parameters and standard errors, but first we'll focus on computing the estimates. This is done using the function `LinearRegression()` in the `sklearn.linear_model` module. The syntax is superficially somewhat similar to using `PCA(...)` in the previous lab (also from an `sklearn` module): we'll configure `LinearRegression(...)` as `slr`, and then fit the model using the `fit(...)` method.

```
In [56]: # configure module
slr = LinearRegression(fit_intercept = False)

# fit model
slr.fit(x_slr, y)
```

```
Out[56]: ▼ LinearRegression
LinearRegression(fit_intercept=False)
```



Notice the argument `fit_intercept = False` -- this configures `slr` *not* to fit an intercept, which is done here since the explanatory variable matrix already includes a column of ones for the intercept term.

## Q1 (a) Retrieve estimated coefficients

After the `fit(...)` method is applied, the coefficient estimates are stored as the `.coef_` attribute of `slr`.

(i) Store the coefficient estimates in `coef` and print them.

```
In [57]: #store the coefficient estimates
coef = slr.coef_

# print coefficient estimates
coef
```

```
Out[57]: array([ 7.51142307, -0.42747211])
```

```
In [58]: grader.check("q1_a_i")
```

```
Out[58]: q1_a_i passed! 🌈
```

The order in which the `.coef_` attribute stores the estimates matches the order of the columns of  $\mathbf{X}$ .

(ii) Which is which?

Identify the intercept and identify the slope. Replace the dots below.

- Intercept estimate:  $\hat{\beta}_0 = \dots$
- Slope estimate:  $\hat{\beta}_1 = \dots$

(Hint: refer back to the plot; only one of the two estimates could possibly be the slope.)

```
In [59]: b0 = 7.51142307
b1 = -0.42747211
```

```
In [60]: grader.check("q1_a_ii")
```

```
Out[60]: q1_a_ii passed! 🙌
```

(iii) Interpret the slope coefficient.

Fill in the blanks below.

Among the countries in the sample, a one-year increase in **the expected years of education for an adult woman** is associated with a **0.427 decrease** in the average number of children per woman.

## Fitted values and residuals

The 'fitted value' for  $y_i$  is the value along the line specified by the model that corresponds to the matching explanatory variable  $x_i$ . In other words:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

These can be obtained by passing the explanatory variable matrix to the `.predict(...)` method. In effect, this returns predictions on the observed data.

```
In [61]: # fitted values
fitted_slr = slr.predict(x_slr)

fitted_slr[0:5]
```

```
Out[61]: array([4.60644126, 1.86804122, 2.33516772, 4.5302748 , 1.99085971])
```

The result is an array with length matching the number of rows in `x_slr`.

```
In [62]: # dimensions
fitted_slr.shape
```

```
Out[62]: (139,)
```

## Q1 (b) Compute the residuals

The 'residuals' -- what's left over -- are the differences between the fitted and observed values of the response:

$$e_i = y_i - \hat{y}_i \quad i\text{th residual}$$

Use `fitted_slr` and `y` to compute an array of the residuals. Store the result as `resid_slr`.

```
In [63]: # residuals
resid_slr = y - fitted_slr

# print
resid_slr
```

Out[63]: Country  
Afghanistan -0.133441  
Albania -0.251041  
Algeria 0.687832  
Angola 0.988725  
Antigua and Barbuda 0.003140  
...  
Uruguay -0.204590  
Uzbekistan 0.224314  
Vanuatu 0.602464  
Zambia 0.859878  
Zimbabwe 0.795048  
Name: fertility\_total, Length: 139, dtype: float64

```
In [64]: grader.check("q1_b")
```

Out[64]: q1\_b passed! 🎉

### Q1 (c) Append fitted values and residuals to the data

In order to keep them handy, add the fitted values and residuals as new columns in `reg_data` with the names `fitted_slr` and `resid_slr`. Print the first few rows of `reg_data` after adding the new variables.

```
In [65]: # append
reg_data['fitted_slr'] = fitted_slr
reg_data['resid_slr'] = resid_slr

# print
reg_data
```

Out[65]:

	fertility_total	educ_expected_yrs_f	hdi	hdi_fac	fitted_slr	resid_slr
Country						
Afghanistan	4.473	6.795722	0.509	1_low	4.606441	-0.133441
Albania	1.617	13.201755	0.792	2_medium	1.868041	-0.251041
Algeria	3.023	12.108990	0.746	2_medium	2.335168	0.687832
Angola	5.519	6.973901	0.582	1_low	4.530275	0.988725
Antigua and Barbuda	1.994	12.914441	0.772	2_medium	1.990860	0.003140
...	...	...	...	...	...	...
Uruguay	1.973	12.477616	0.816	2_medium	2.177590	-0.204590
Uzbekistan	2.604	12.004848	0.717	2_medium	2.379686	0.224314
Vanuatu	3.782	10.133731	0.603	1_low	3.179536	0.602464
Zambia	4.633	8.745135	0.582	1_low	3.773122	0.859878
Zimbabwe	3.615	10.974918	0.569	1_low	2.819952	0.795048

139 rows x 6 columns

```
In [66]: grader.check("q1_c")
```

Out[66]: q1\_c passed! 🙌

There's one more estimate to compute: the error variance,  $\sigma^2$ ! This can be estimated based on the residual variance:

$$\hat{\sigma}^2 = \frac{1}{n - p} \sum_{i=1}^n e_i^2 = \frac{n - 1}{n - p} S_e^2$$

In this expression,  $e_i$  are the residuals,  $n$  and  $p$  are the dimensions of  $\mathbf{X}$ , and  $S_e^2$  is the residual (sample) variance.

```
In [67]: # store n (number of observations) and p (number of betas)
n, p = x_slr.shape

# compute estimate of error variance
sigma2_hat = ((n - 1)/(n - p)) * resid_slr.var()
```

### Uncertainty quantification

It was noted in lecture that the variances and covariances of  $\hat{\beta}_0, \hat{\beta}_1$  are given by the matrix:



$$\sigma^2(\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} \text{var}\hat{\beta}_0 & \text{cov}\left(\hat{\beta}_0, \hat{\beta}_1\right) \\ \text{cov}\left(\hat{\beta}_1, \hat{\beta}_0\right) & \text{var}\hat{\beta}_1 \end{bmatrix}$$

So we can *estimate* these quantities, which quantify the variation and covariation of the estimated coefficients, by plugging in the estimated error variance and computing:

$$\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$$

```
In [68]: # compute X'X
xtx_slr = x_slr.transpose().dot(x_slr)

# compute matrix of parameter variances/covariances: estimated error variance x (X'X)^{-1}
slrcoef_vcov = np.linalg.inv(xtx_slr) * sigma2_hat

# print
slrcoef_vcov
```

```
Out[68]: array([[ 0.0688336 , -0.0057818 ],
               [-0.0057818 ,  0.00050894]])
```

### Q1 (d) Standard errors

The *standard errors* for the estimates  $\hat{\beta}_0, \hat{\beta}_1$  are the square roots of their estimated variances. Compute these by retrieving the diagonal elements of `slrcoef_vcov` (the estimated variances) and taking square roots.

Store the result (an array with two values) as `slrcoef_se` and print the array.

(*Hint*: use `.diagonal()` from numpy.)

```
In [69]: # take square root of matrix diagonals to get standard errors
slrcoef_se = np.sqrt(slrcoef_vcov.diagonal())

# print
slrcoef_se
```

```
Out[69]: array([0.26236158, 0.02255976])
```

```
In [70]: grader.check("q1_d")
```

```
Out[70]: q1_d passed! 🌟
```

We can now report the results of model fitting in an organized fasion: we'll make a dataframe `slrcoef_table` with two columns, 'coefficient estimate' and 'standard error', and index the rows of the dataframe by the coefficient names (intercept and education).

```
In [71]: # create table
slrcoef_table = pd.DataFrame(
    data = {'coefficient estimate': slr.coef_, 'standard error': slrcoef_se},
    index= ['intercept', 'education']
)

# print
slrcoef_table
```

```
Out[71]:
```

	coefficient estimate	standard error
intercept	7.511423	0.262362
education	-0.427472	0.022560

Lastly, a standard metric often reported with linear models is the  $R^2$  score, which is interpreted as the proportion of variation in the response captured by the model. It is straightforward to compute using the `r2_score(...)` function from `sklearn.metrics`:

```
In [72]: # compute R-squared
r2_score(reg_data.fertility_total, reg_data.fitted_slr)
```

```
Out[72]: 0.723814308777671
```

But just so you have a sense of what it is, a direct calculation is shown below:

```
In [73]: # R-squared 'by hand'
(y.var() - resid_slr.var())/y.var()
```

```
Out[73]: 0.7238143087776711
```

The metric is simply the difference between the raw variation in the response and residual variation, as a proportion of the variation in the response. If the model fits well, the residual variation will be small, in which case this proportion will be closer to 1.

## Q1 (e) Visualize!

Now that you've reported the numerical results of model fitting, have a direct look at the results relative to the data scatter.

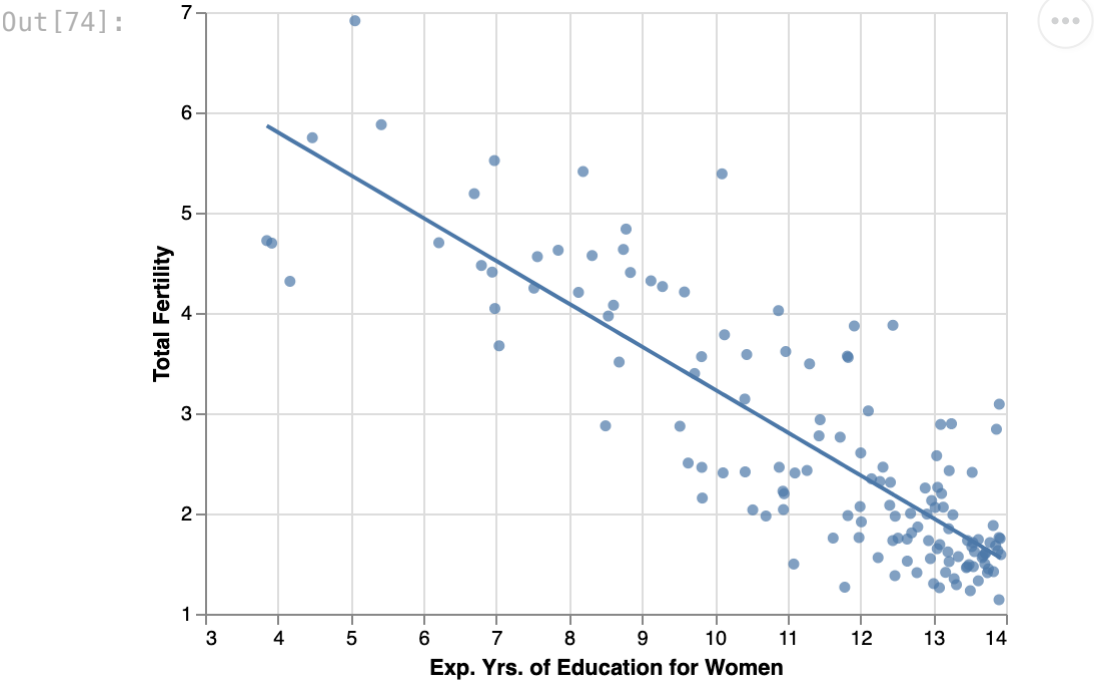
(i) Construct a line plot of the fitted values against expected years of education.

Remember that the fitted values are stored as `fitted_slr` in the dataframe `reg_data` containing the actual data.

Construct the line as `slr_line` and layer this on top of your plot `simple_scatter` from part 0 above.

```
In [74]: # construct line plot
slr_line = alt.Chart(reg_data).mark_line().encode(
    x = alt.X('educ_expected_yrs_f'),
    y = alt.Y('fitted_slr'))

# layer
slr_line + simple_scatter
```



(ii) Comment on the plot.

How well does the model seem to describe the data overall?

The model seems to fit the data fairly well, as it captures the overall trend and picture of the data, but I do believe that there are other models that would probably to even better.

---

## 2. Multiple linear regression

Now let's consider adding the human development factor -- a categorical variable -- to the model. The factor column `hdi_fac` can't be input directly as a variable into the explanatory variable matrix  $\mathbf{X}$ , because it doesn't make sense to multiply a coefficient (number) by a factor label (character string), so we'll need to do a little preprocessing to incorporate the factor into a multiple regression model.

### Encoding categorical variables

Categorical variables must be *encoded* in the explanatory variable matrix by *indicator functions*; each level of the category will have a corresponding coefficient. In this case, the factor (HDI) had three levels (low, medium, and high), so we'll construct a matrix  $\mathbf{X}$  with entries

$$\begin{bmatrix} 1 & x_1 & I(hdi_1 = med) & I(hdi_1 = high) \\ 1 & x_2 & I(hdi_2 = med) & I(hdi_2 = high) \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & I(hdi_n = med) & I(hdi_n = high) \end{bmatrix}$$

where

- $I(hdi_i = med)$  is 1 if the  $i$ th HDI level is medium, and 0 otherwise and
- $I(hdi_i = high)$  is 1 if the  $i$ th HDI level is high, and 0 otherwise.

The low level of the HDI factor is not represented explicitly, but if  $I(hdi_i = med) = I(hdi_i = high) = 0$ , then  $hdi_i = low$ ; in other words, if the  $i$ th level of the HDI factor is neither medium nor high, then it must be low.

It takes a little head-scratching to figure out what these indicators do, but **their effect in the model is to allow the intercept to change depending on the level of the factor.**

These indicators functions are sometimes called 'dummy' variables. They're simple to obtain using `pd.get_dummies(...)`:

```
In [75]: # encode hdi factor by indicator variables
hdi_df = pd.get_dummies(reg_data.hdi_fac, drop_first = True)
```

```
# preview
hdi_df.head(6)
```

Out[75]:

	2_medium	3_high
Country		
Afghanistan	0	0
Albania	1	0
Algeria	1	0
Angola	0	0
Antigua and Barbuda	1	0
Argentina	0	1

Now compare this with the underlying factor:

```
In [76]: reg_data[['hdi_fac']].head(6)
```

Out[76]:

	hdi_fac
Country	
Afghanistan	1_low
Albania	2_medium
Algeria	2_medium
Angola	1_low
Antigua and Barbuda	2_medium
Argentina	3_high

## Interaction terms

We'll also add *interaction* terms -- a fancy phrase for products of explanatory variables -- between the indicators and the education variable. **The effect of interactions is to allow the slope to change depending on the level of the factor.**

The cell below computes these interaction terms and appends them, along with the indicators, to the education variable.

```
In [77]: # compute interaction terms
interaction_df = hdi_df.multiply(reg_data.educ_expected_yrs_f, axis = 0)
interaction_df.columns = ['2_medium_x_educ', '3_high_x_educ']

# append indicators to data
x_mlr_df = pd.concat([reg_data.educ_expected_yrs_f, hdi_df, interaction_df], ignore_index = False, axis = 1)

# preview
x_mlr_df.head(4)
```

Out[77]:

	educ_expected_yrs_f	2_medium	3_high	2_medium_x_educ	3_high_x_educ
Country					
Afghanistan	6.795722	0	0	0.000000	0.0
Albania	13.201755	1	0	13.201755	0.0
Algeria	12.108990	1	0	12.108990	0.0
Angola	6.973901	0	0	0.000000	0.0

Notice that those interaction terms are simply the elementwise products of each indicator column with the education variable column.

## Model fitting

With the encoded categorical variable for HDI, and the interaction terms, the linear model is:

$$\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} 1 & x_1 & I(hdi_1 = med) & I(hdi_1 = high) & x_1 I(hdi_1 = med) & x_1 I(hdi_1 = high) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & I(hdi_n = med) & I(hdi_n = high) & x_n I(hdi_n = med) & x_n I(hdi_n = high) \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \end{bmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}}_{\boldsymbol{\epsilon}}$$

Which is matrix form for:

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 I(hdi_1 = med) + \beta_3 I(hdi_1 = high) + \beta_4 x_1 I(hdi_1 = med) + \beta_5 x_1 I(hdi_1 = high) + \epsilon_i \quad i = 1, \dots, n$$

While seemingly more complicated, and perhaps a little trickier to interpret mathematically, the syntactical pattern for fitting the model is exactly the same as the simple linear model.

```
In [78]: # add intercept column
x_mlr = add_dummy_feature(x_mlr_df, value = 1)

# configure module
mlr = LinearRegression(fit_intercept = False)

# fit model
mlr.fit(x_mlr, y)
```

```
Out[78]: LinearRegression
LinearRegression(fit_intercept=False)
```

Q2 (a) Summarize model fit

Follow the examples from the simple linear model to carry out the following:

- 1. Compute the error variance estimate  $\hat{\sigma}^2$  and store it in `sigma2_hat`.
- 2. Compute the standard errors for each coefficient and store in `mlrcoef_se`.
- 3. Construct a dataframe called `mlrcoef_table` showing both the estimates and standard errors.

(Suggestion: copy the codes from part 1 into a single cell below for reference.)

```
In [79]: # store dimensions of explanatory variable matrix
n, p = x_mlr.shape

# compute residual variance
fitted_mlr = mlr.predict(x_mlr)
resid_mlr = y - fitted_mlr
sigma2_hat = ((n - 1)/(n - p)) * resid_mlr.var()

# compute standard errors
xtx_mlr = x_mlr.transpose().dot(x_mlr)

mlrcoef_vcov = np.linalg.inv(xtx_mlr) * sigma2_hat
mlrcoef_se = np.sqrt(mlrcoef_vcov.diagonal())

# construct coefficient table
mlrcoef_table = pd.DataFrame(
    data = {'coefficient estimate': mlr.coef_, 'standard error': mlrcoef_se},
    index = ['Intercept', 'Education', 'Medium', 'High', 'Medium x Education', 'High x Education']
)
# print
mlrcoef_table
```

Out[79]:

	coefficient estimate	standard error
Intercept	7.049013	0.369182
Education	-0.351265	0.041062
Medium	-2.057722	0.898528
High	-4.504779	2.485328
Medium x Education	0.125183	0.079162
High x Education	0.286587	0.187939

```
In [80]: grader.check("q2_a")
```

```
Out[80]: q2_a passed! ✨
```

Q2 (b) Fit metric

Now let's compare the fit of the MLR model to that of the SLR model.

(i) Calculate the  $R^2$  score for the multiple linear regression model and store it in `R_2`.

```
In [81]: # solution
R_2 = (y.var() - resid_mlr.var())/y.var()

R_2
```

```
Out[81]: 0.7657210526482656
```

```
In [82]: grader.check("q2_b_i")
```

```
Out[82]: q2_b_i passed! 🎉
```

(ii) Does the fit seem to have improved appreciably?

Answer

Yes, it does since the R-squared value has increased.

## Q2 (c) Visualize

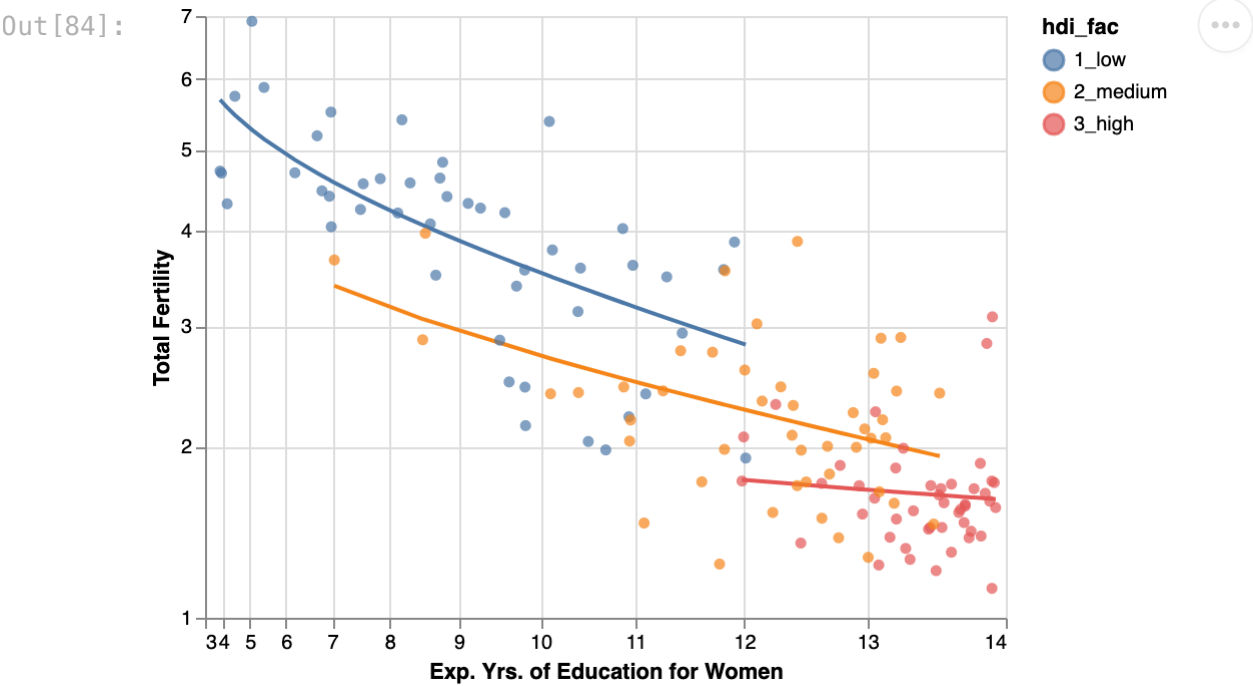
Now visualize the model outputs.

1. First append the fitted values (if you didn't above) as a new column named `fitted_mlr` to `reg_data`.
2. Construct a line plot of the fitted values (y axis) against education (x axis) and map the HDI factor to the color aesthetic. This should produce a plot with three lines, one for each level of HDI. Store the result as `mlr_lines`.
3. Layer `mlr_lines` on top of `scatter_scaled`.

```
In [84]: # append fitted values to original data
reg_data['fitted_mlr'] = fitted_mlr

# construct line plot
mlr_lines = alt.Chart(reg_data).mark_line().encode(
    x = alt.X('educ_expected_yrs_f'),
    y = alt.Y('fitted_mlr'),
    color = 'hdi_fac:N'
)

# layer
mlr_lines + scatter_scaled
```



## Q2 (d) Interpret your graphic

What does the multiple linear model suggest about the relationship between fertility rate and expected years of education for high-HDI countries as compared with low-HDI countries?

(Hint: focus on how the slopes of the two lines compare.)

### Answer

According to the multiple linear model for high-HDI and low-HDI countries as fitted above, we are able to see a slightly different relationship between the two when it comes to their fertility rates. The blue line, or the low-HDI countries seem to have a much steeper declining line, indicating that there is a much more negative relationship in those countries when compared to the red line which corresponds to the high-HDI countries. For them, a few extra years of education has a big impact on the total fertility rate. The bottom red line is much more flat and does not decline as much as the blue one, meaning that more education for women in high countries does not have as drastic of a change in the total fertility for those individuals.

## Comment

These data are definitely *not* a representative sample of any particular population of nations -- the countries (observational units) are conveniently chosen based on which countries reported data. So there is no scope of inference here.

Although we can't claim that, for example, 'the mean fertility rate decreases with education at a rate of -0.35 children per woman per expected year of education in low-HDI countries', we *can* say '*among the countries reporting data*, the mean fertility rate decreases with education at a rate of -0.35 children per woman per expected year of education in low-HDI countries'. This is a nice example of how a model might be used in a descriptive capacity.

---

## Submission Checklist

1. Save file to confirm all changes are on disk
2. Run *Kernel > Restart & Run All* to execute all code from top to bottom
3. Save file again to write any new output to disk
4. Select *File > Download as > HTML*.
5. Open in Google Chrome and print to PDF on A3 paper in portrait orientation.

To double-check your work, the cell below will rerun all of the autograder tests.

```
In [85]: grader.check_all()
```

```
Out[85]: q1_a_i results: All test cases passed!
```

```
q1_a_ii results: All test cases passed!
```

```
q1_b results: All test cases passed!
```

```
q1_c results: All test cases passed!
```

```
q1_d results: All test cases passed!
```

```
q2_a results: All test cases passed!
```

```
q2_b_i results: All test cases passed!
```