



The general problem

- We have n data points, $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)$, where each \vec{x}_i is a feature vector of d features:

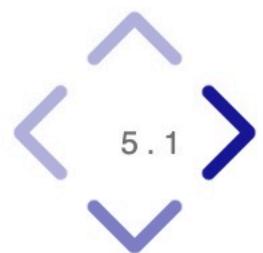
$$\vec{x}_i = \begin{bmatrix} x_i^{(1)} \\ x_i^{(2)} \\ \vdots \\ x_i^{(d)} \end{bmatrix}$$

e.g. predicting commute time given

- departure hour

- day of month

$$\vec{x}_5 = \begin{bmatrix} 8.5 \\ 13. \end{bmatrix}$$





- We want to find a good linear hypothesis function:

$$H(\vec{x}_i) = w_0 + w_1 x_i^{(1)} + w_2 x_i^{(2)} + \dots + w_d x_i^{(d)} = \vec{w} \cdot \text{Aug}(\vec{x}_i)$$

- Specifically, we want to find the optimal parameters, $w_0^*, w_1^*, \dots, w_d^*$ that minimize mean squared error:

mean squared error

$$\begin{aligned}
 R_{\text{sq}}(\vec{w}) &= \frac{1}{n} \sum_{i=1}^n (y_i - H(\vec{x}_i))^2 \\
 &= \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i^{(1)} + w_2 x_i^{(2)} + \dots + w_d x_i^{(d)}))^2 \\
 &= \frac{1}{n} \sum_{i=1}^n (y_i - \text{Aug}(\vec{x}_i) \cdot \vec{w})^2 \\
 &= \frac{1}{n} \|\vec{y} - X\vec{w}\|^2
 \end{aligned}$$

X : "design matrix"





The general solution

- Define the **design matrix** $X \in \mathbb{R}^{n \times (d+1)}$ and **observation vector** $\vec{y} \in \mathbb{R}^n$:

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(d)} \\ 1 & x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(d)} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(d)} \end{bmatrix} = \begin{bmatrix} \text{Aug}(\vec{x}_1)^T \\ \text{Aug}(\vec{x}_2)^T \\ \vdots \\ \text{Aug}(\vec{x}_n)^T \end{bmatrix}$$

n × d + 1

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

\vec{h} , which contains predictions,
 $\vec{X}\vec{w}$, which is a linear combination
of X 's columns.

$$\begin{bmatrix} 1 & \vec{x}_n^{(1)} & \vec{x}_n^{(2)} & \dots & \vec{x}_n^{(d)} \end{bmatrix} \quad \text{Aug}(\vec{x}_n)^T$$

- Then, solve the **normal equations** to find the optimal parameter vector, \vec{w}^* :

solution : orthogonal projection
of \vec{y} onto $\text{span}(X)$

observe: sum of each row is always exactly important

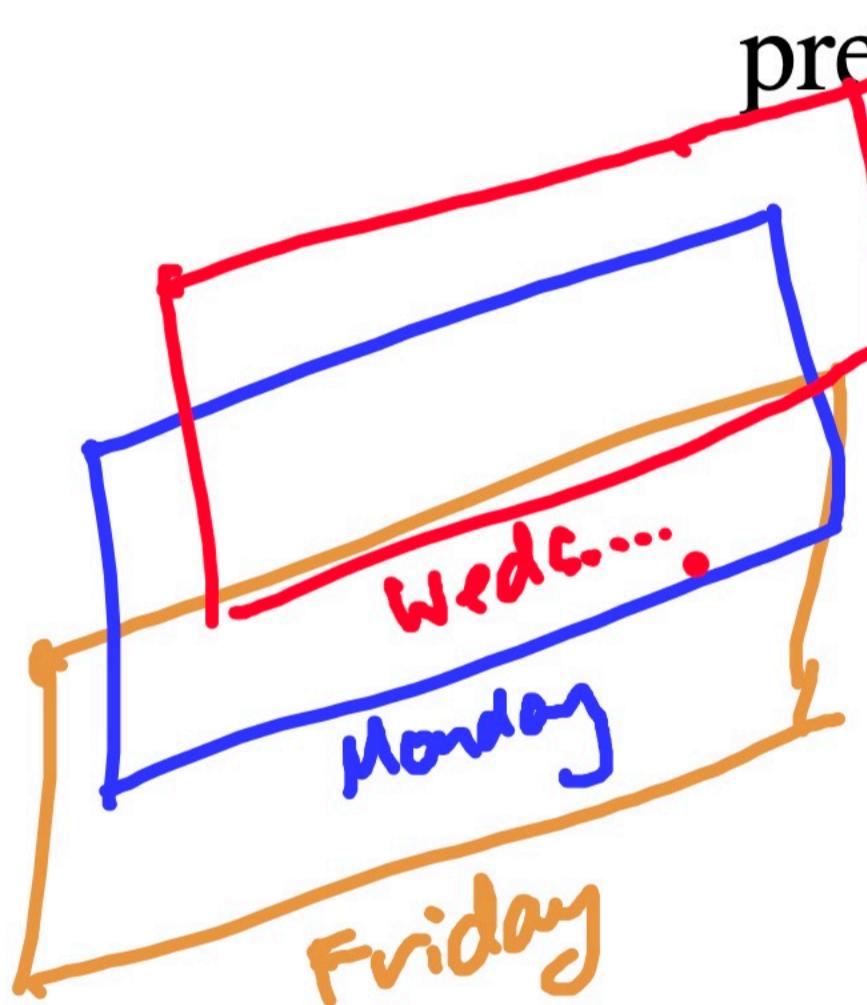
day	day == Mon	day == Tue	day == Wed	day == Thu	day == Fri
Mon	1	0	0	0	0
Tue	0	1	0	0	0
Mon	1	0	0	0	0
...
Mon	1	0	0	0	0
Tue	0	1	0	0	0
Thu	0	0	0	1	0



Visualizing our latest model

- Our trained linear model to predict commute time given 'departure_hour', 'day_of_month', and 'day' (Mon, Tue, Wed, or Thu) is:

$$\begin{aligned} \text{pred. commute time}_i &= 134 - 8.42 \cdot \text{departure hour}_i \\ &\quad - 0.03 \cdot \text{day of month}_i \\ &\quad + 5.09 \cdot \text{day}_i == \text{Mon} \\ &\quad + 16.38 \cdot \text{day}_i == \text{Tue} \\ &\quad + 5.12 \cdot \text{day}_i == \text{Wed} \\ &\quad + 11.5 \cdot \text{day}_i == \text{Thu} \end{aligned}$$



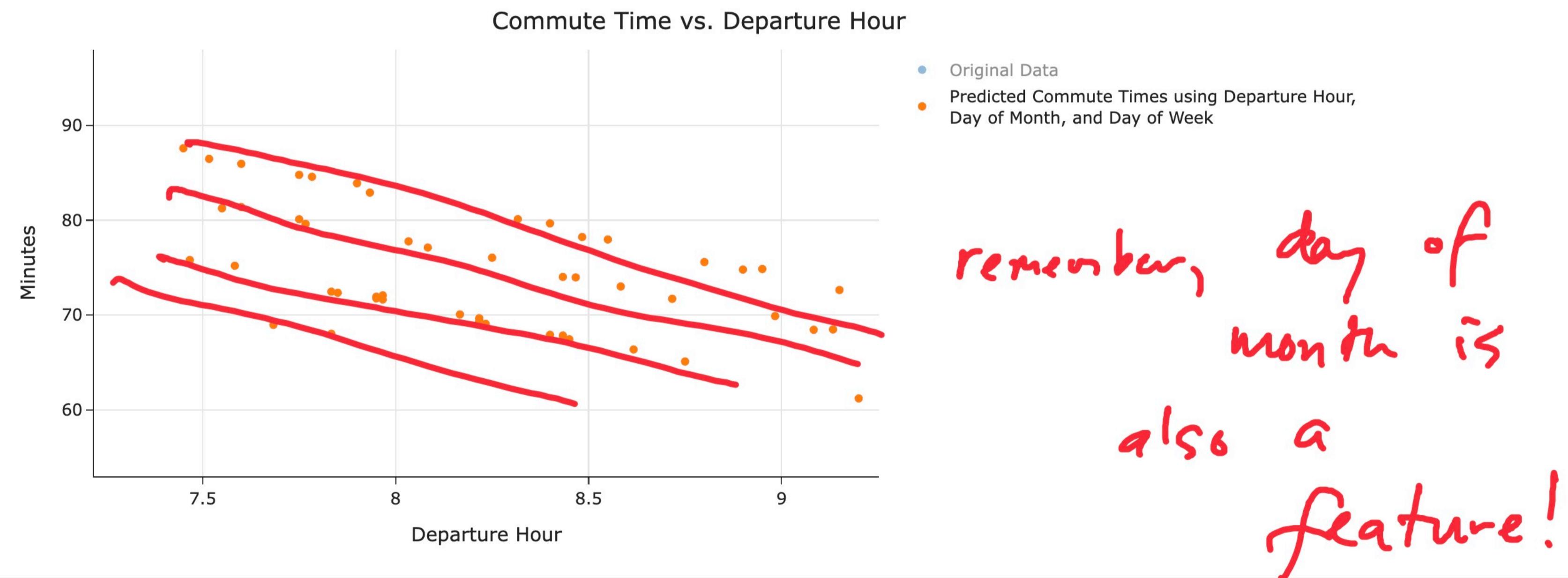
- Since we have 6 features here, we'd need 7 dimensions to graph our model.
- But, as we see in Homework 7, Question 5, our model is really a collection of **five parallel planes** in 3D, all with slightly different z -intercepts!





- If we want to visualize in 2D, we need to pick a single feature to place on the x -axis.

```
In [16]: 1 fig = go.Figure()
2 fig.add_trace(go.Scatter(x=df['departure_hour'], y=df['minutes'],
3                           mode='markers', name='Original Data'))
4 fig.add_trace(go.Scatter(x=df['departure_hour'], y=model_with_ohe.predict(X_for_ohe),
5                           mode='markers', name='Predicted Commute Times using Departure Hour, <br>Day of Month, and Day of Week'))
6 fig.update_layout(showlegend=True, title='Commute Time vs. Departure Hour',
7                   xaxis_title='Departure Hour', yaxis_title='Minutes', width=1000)
```



- Despite being a linear model, why doesn't this model look like a straight line?

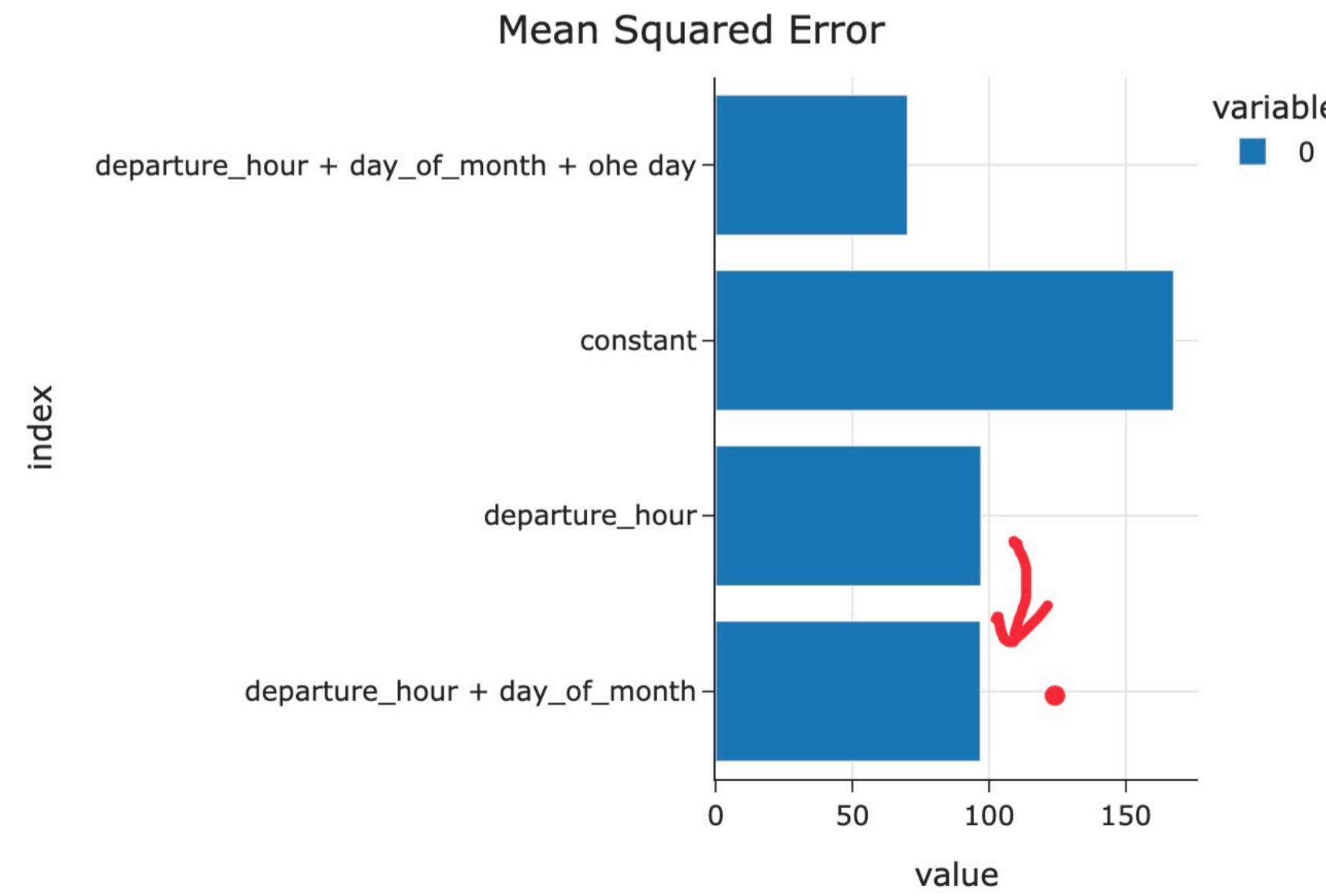




Comparing our latest model to earlier models

- Let's see how the inclusion of the day of the week impacts the quality of our predictions.

```
In [21]: 1 mse_dict['departure_hour + day_of_month + ohe day'] = mean_squared_error(  
2     df['minutes'],  
3     model_with_ohe.predict(X_for_ohe)  
4 )  
5 pd.Series(mse_dict).plot(kind='barh', title='Mean Squared Error')
```



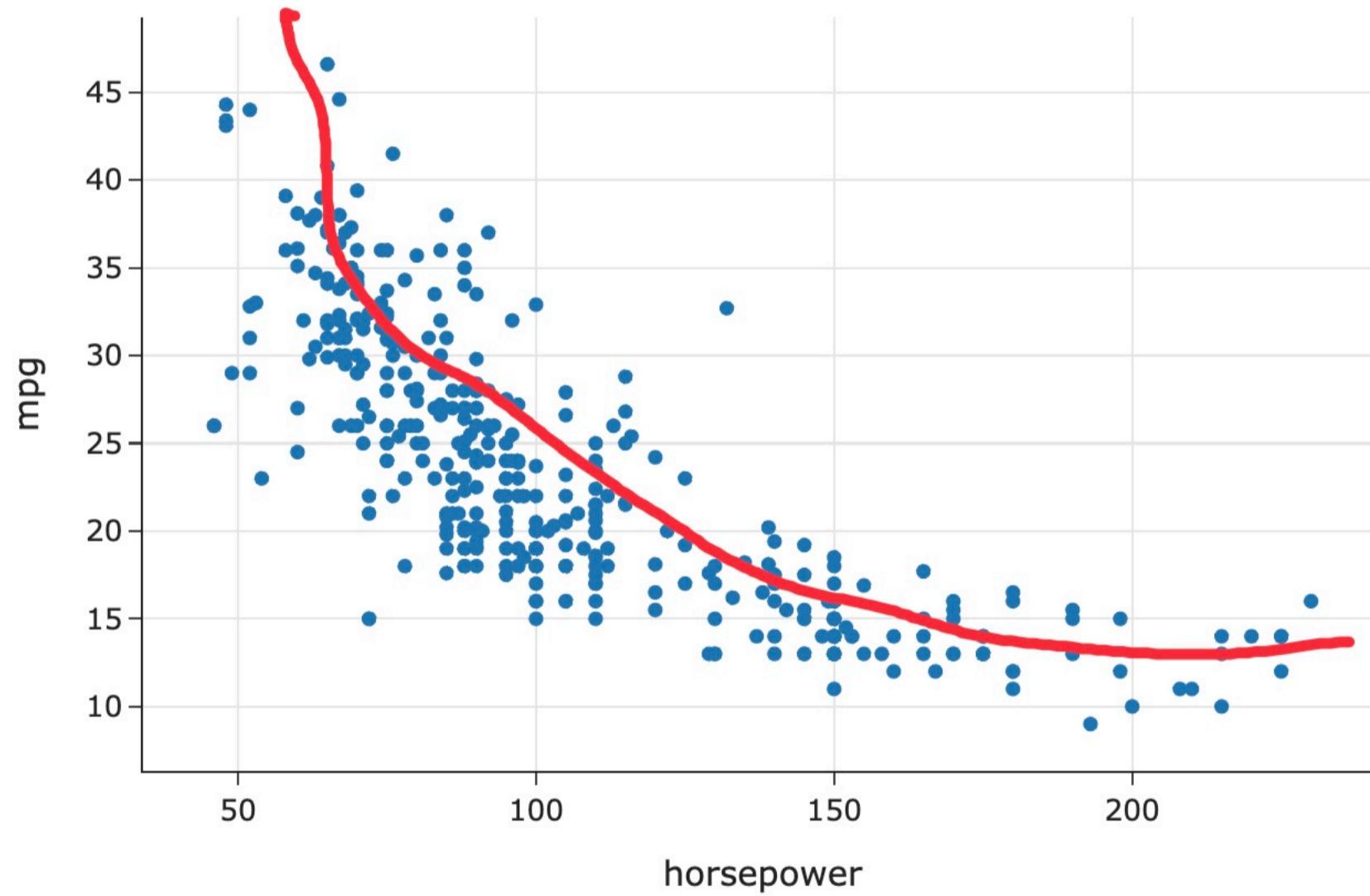
day of month
kind of useless,
but
day of week is
very useful!





The relationship between 'horsepower' and 'mpg'

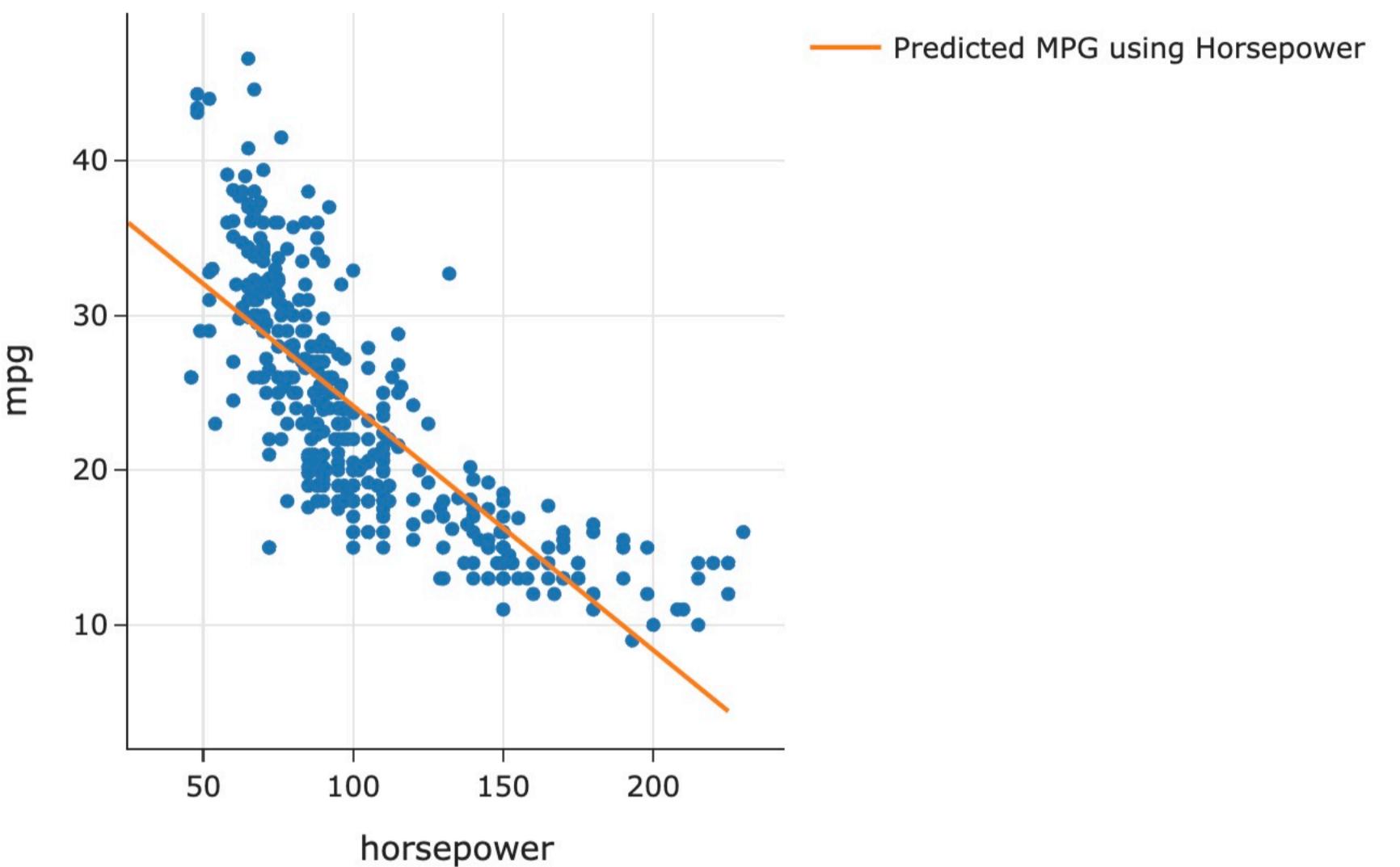
```
In [25]: 1 px.scatter(mpg, x='horsepower', y='mpg')
```



decay.



```
4 x=hp_points['horsepower'],
5 y=car_model.predict(hp_points),
6 mode='lines',
7 name='Predicted MPG using Horsepower'
8 ))
```

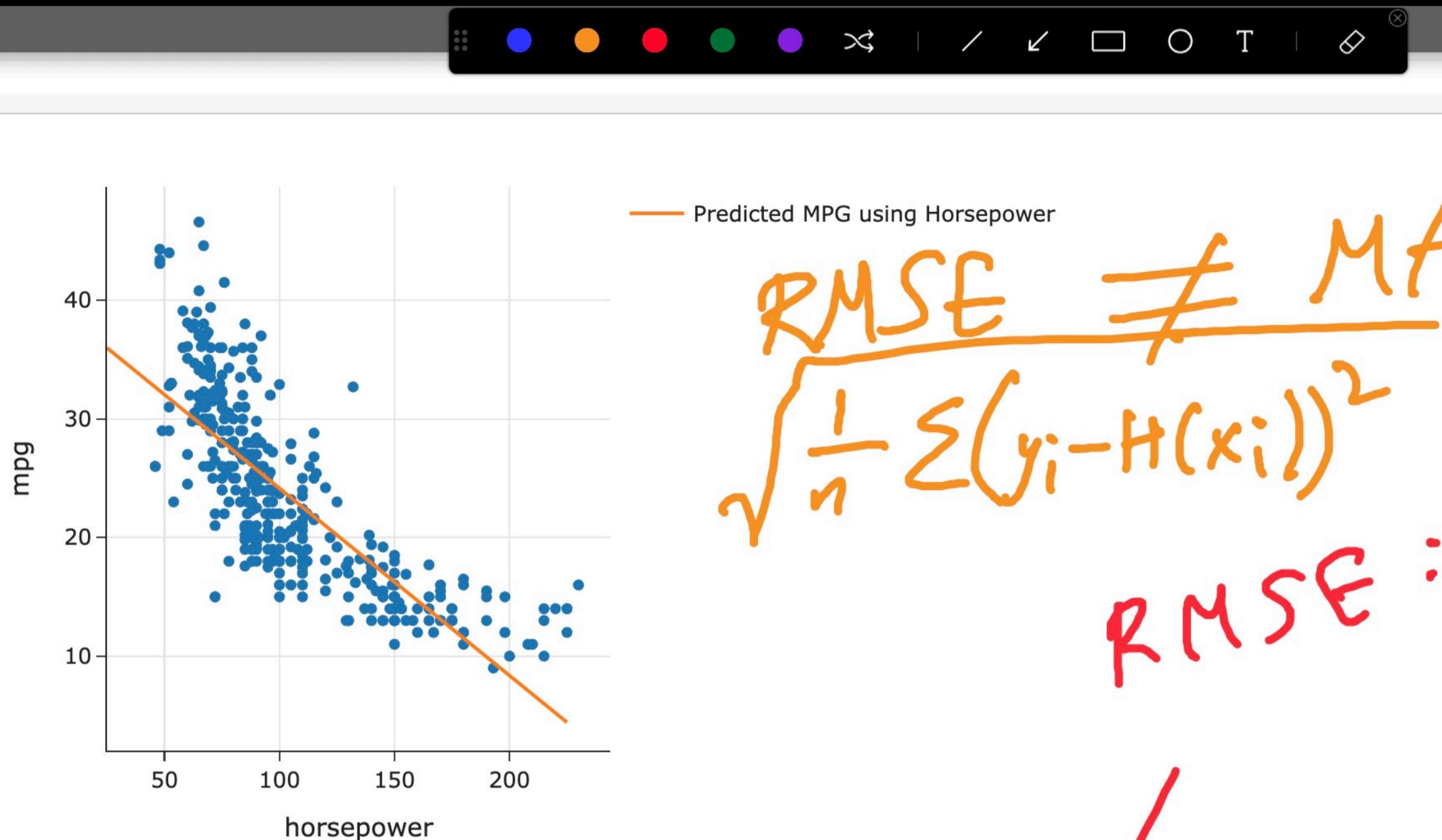


units: mpg²

- Our regression line doesn't capture the curvature in the relationship between 'horsepower' and 'mpg'.

```
In [28]: 1 # As a baseline:  
2 mean_squared_error(mpg['mpg'], car_model.predict(mpg[['horsepower']]))
```

Out [28]: 23.943662938603108



- Our regression line doesn't capture the curvature in the relationship between 'horsepower' and 'mpg'.

```
In [29]: 1 # root mean squared error  
2 mean_squared_error(mpg['mpg'], car_model.predict(mpg[['horsepower']])) * 0.5
```

Out [29]: 4.893226230065713

```
In [28]: 1 # As a baseline:  
2 mean_squared_error(mpg['mpg'], car_model.predict(mpg[['horsepower']]))
```

Out [28]: 23.943662938603108

Linear in the parameters

- Using linear regression, we can fit hypothesis functions like

$$H(x_i) = w_0 + w_1 x_i + w_2 x_i^2$$

~~$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$~~

ysis functions like:

$$H(\vec{x}_i) = w_1 e^{-x_i^{(1)^2}} + w_2 \cos(x_i^{(2)} + \pi) + w_3 \frac{\log 2x_i^{(3)}}{x_i^{(2)}}$$

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$

This includes all polynomials, for example. These are all **linear combinations of (just) features**

- For any of the above examples, we **could** express our model as $\vec{w} \cdot \text{Aug}(\vec{x}_i)$, for some carefully chosen feature vector \vec{x}_i ,
and that's all that `LinearRegression` in `sklearn` needs.

We can be creative

What we put in the `X` argument to `model.fit` is up to us.

We can be creative
with \vec{x}_i

- Using linear regression, we **can't** fit hypothesis functions like:

$$H(x_i) = w_0 + e^{w_1 x_i}$$

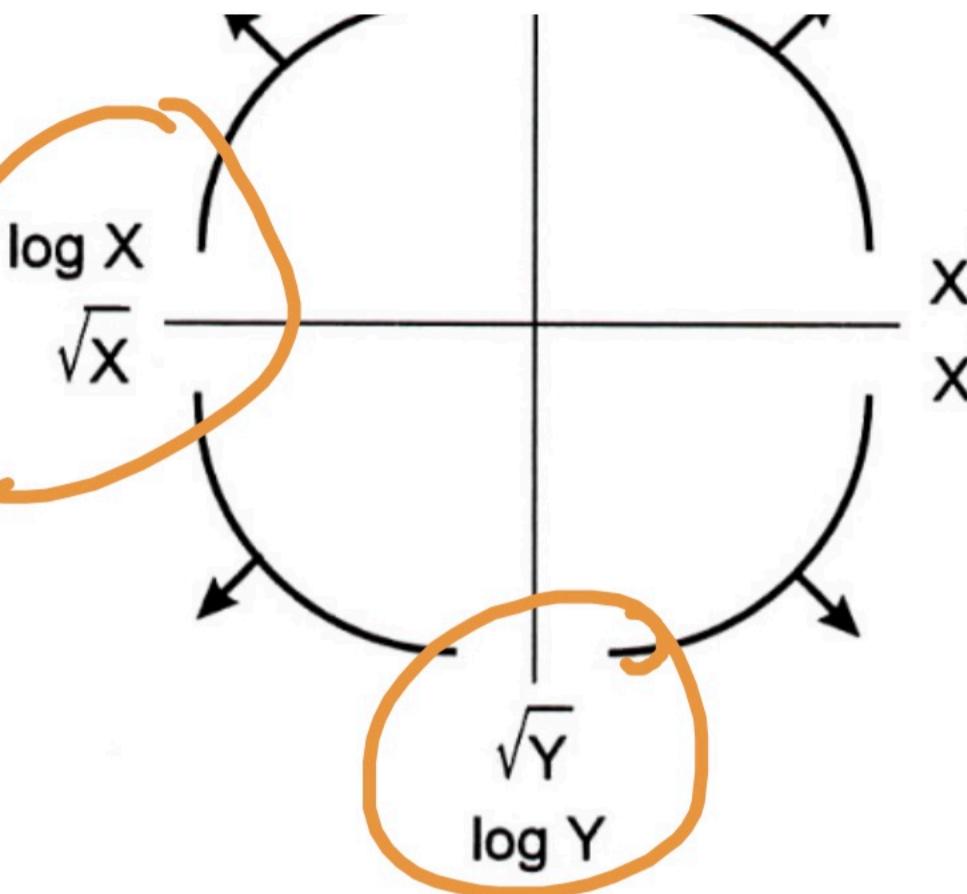
$$H(\vec{x}_i) = w_0 + \sin(w_1 x_i^{(1)} + w_2 x_i^{(2)})$$

doesn't work

These are **not** linear combinations of just features

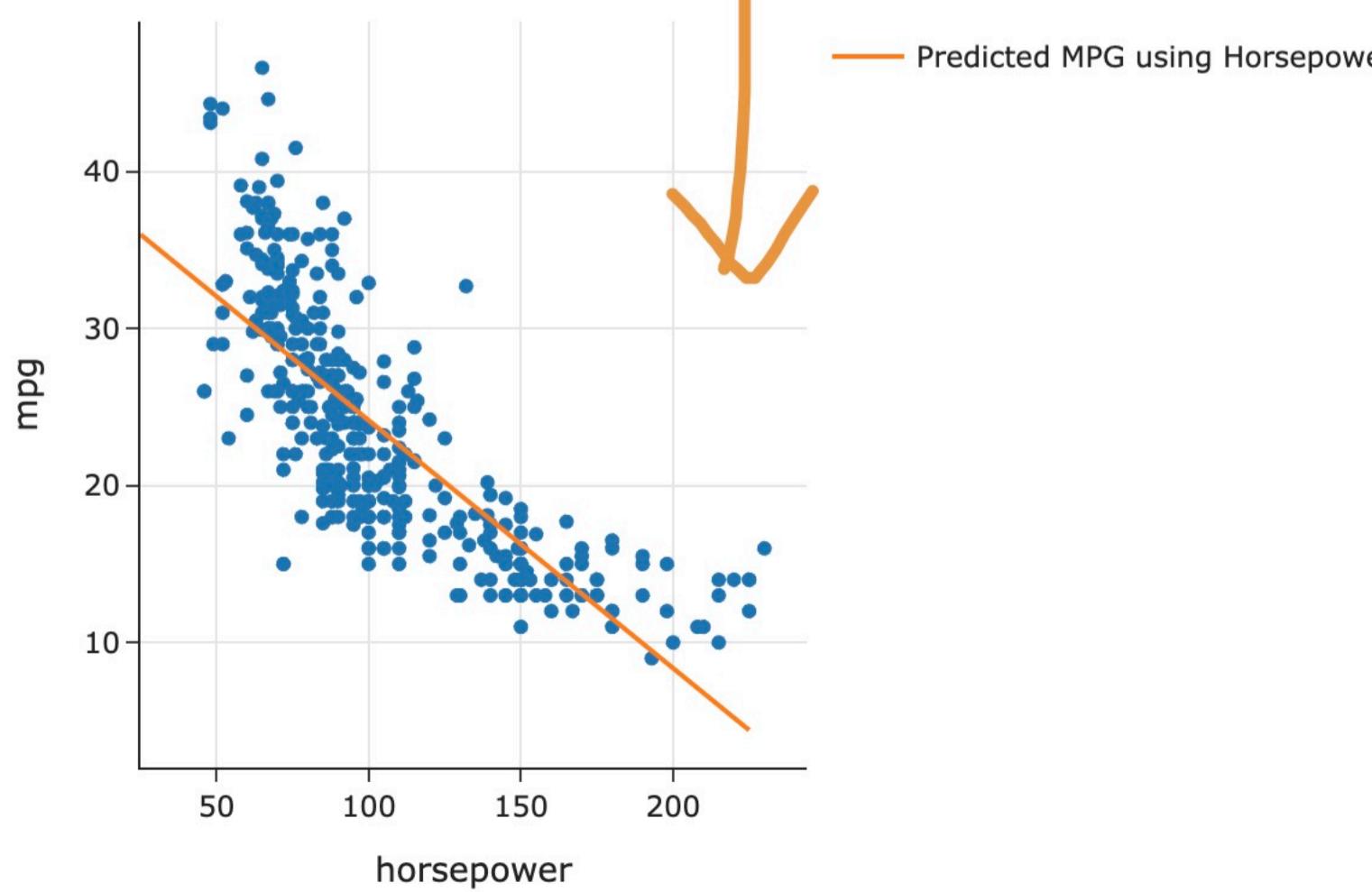


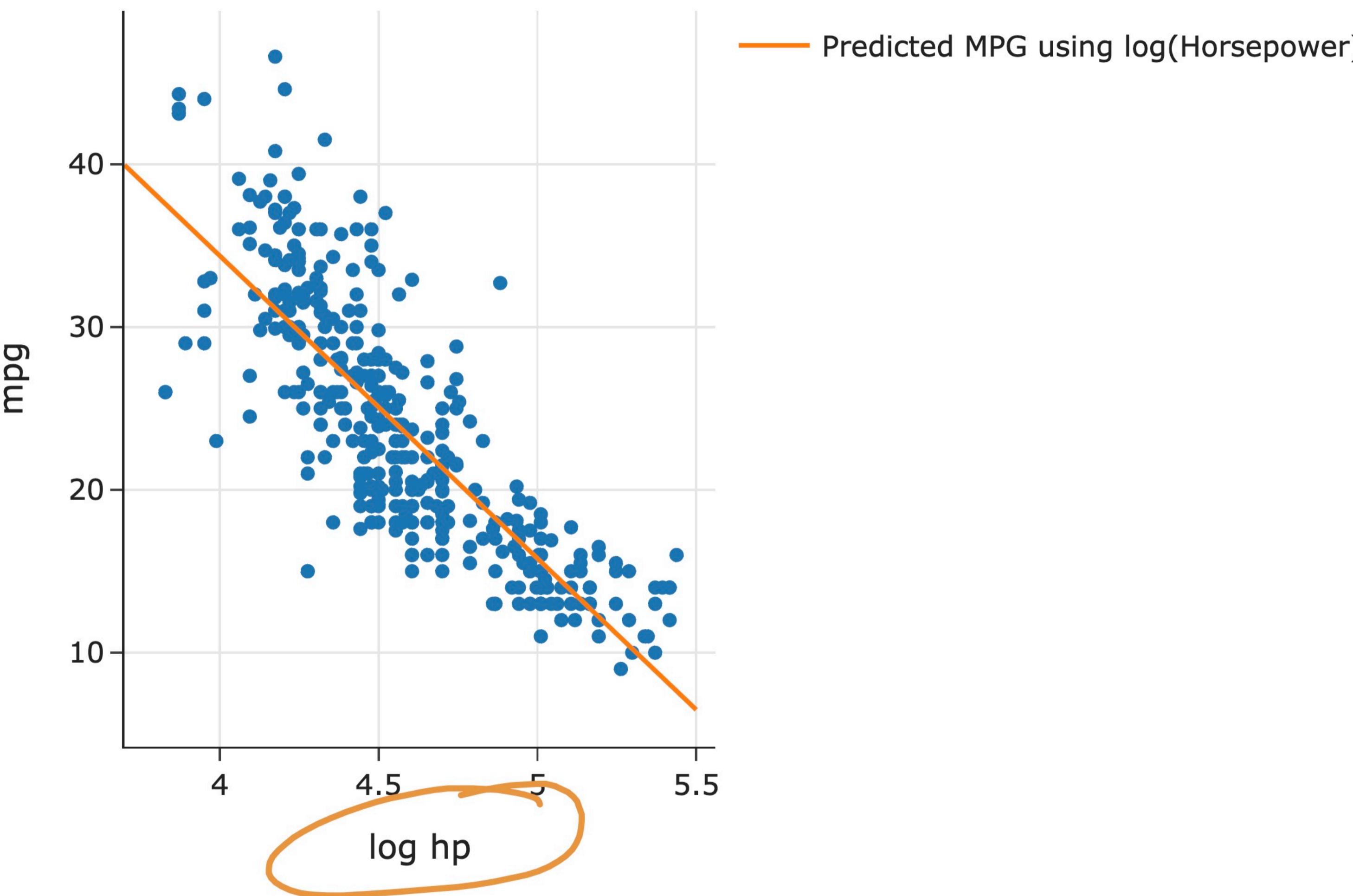
telling me
which features
to create.

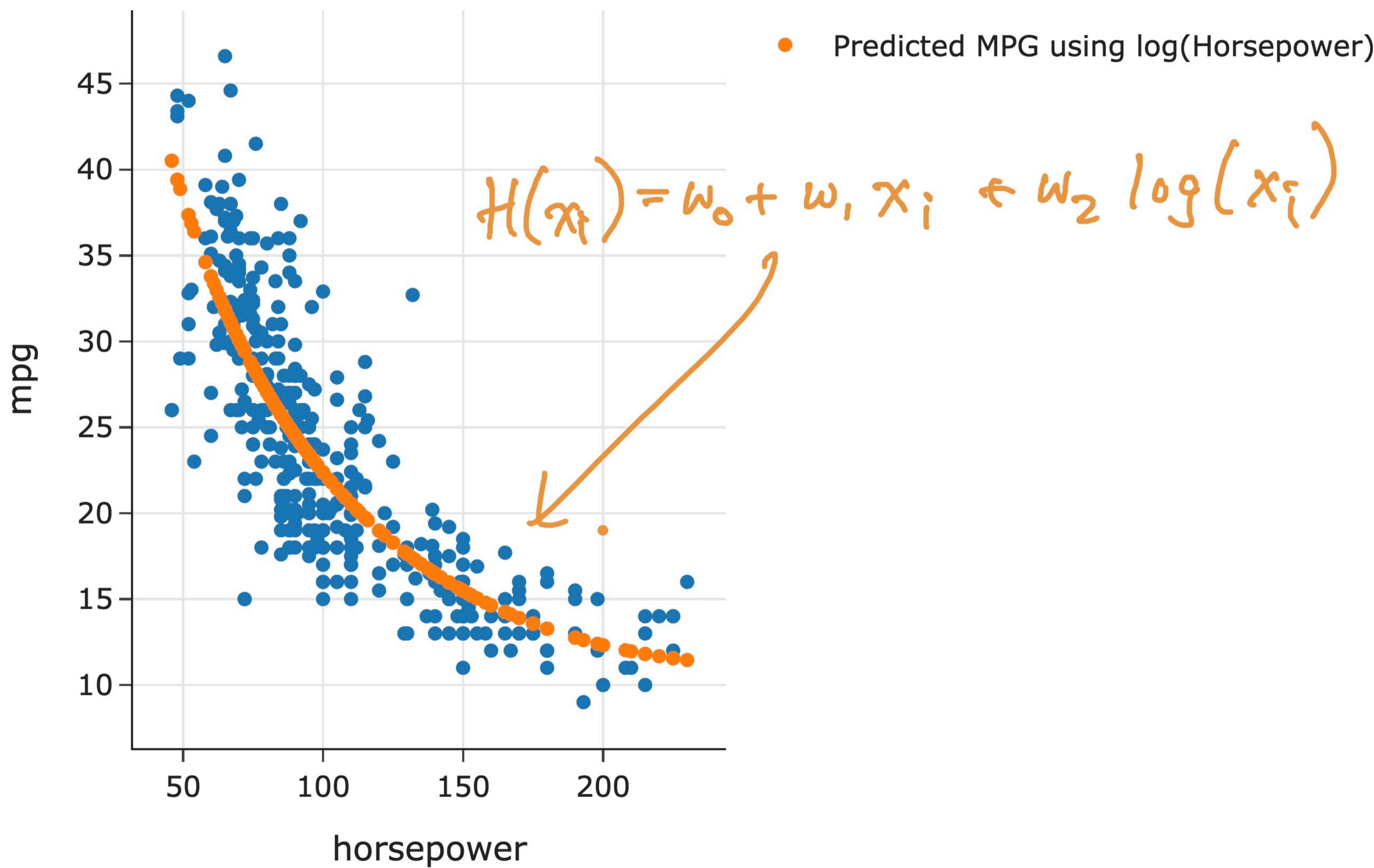


- Why? We're working with linear models. The more linear our data looks in terms of its features, the better we'll be able to model the data.

In [30]: 1 fig









Question 🤔 (Answer at practicaldsc.org/q)

Which hypothesis function is **not** linear in the parameters?

- A. $H(\vec{x}_i) = w_1(x_i^{(1)}x_i^{(2)}) + \frac{w_2}{x_i^{(1)}} \sin(x_i^{(2)})$

- B. $H(\vec{x}_i) = 2^{\omega_1} x_i^{(1)}$

$$\beta_1 = 2^{\omega_1}, \quad H(x_i) = \beta_1 x_i^{(1)}$$

transformation.

- C. $H(\vec{x}_i) = \vec{w} \cdot \text{Aug}(\vec{x}_i)$

- D. $H(\vec{x}_i) = w_1 \cos(x_i^{(1)}) + w_2 2^{x_i^{(2)}} \log x_i^{(3)}$

- E. More than one of the above.

$$\sum w_j \square$$



$$\hat{P} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{z}$$

β_0, β_1

$$\vec{z} = \begin{bmatrix} \log y_1 \\ \log y_2 \\ \vdots \\ \log y_n \end{bmatrix}$$

$$\vec{x} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

$$H(x_i) = w_0 e^{w_1 x_i}$$

log of BS

$$\log H(x_i) = \log(w_0) + w_1 x_i$$

define:

$$z_i = \log y_i \quad T(x_i) = \log H(x_i)$$
$$\beta_0 = \log w_0$$
$$\beta_1 = w_1$$

Then:

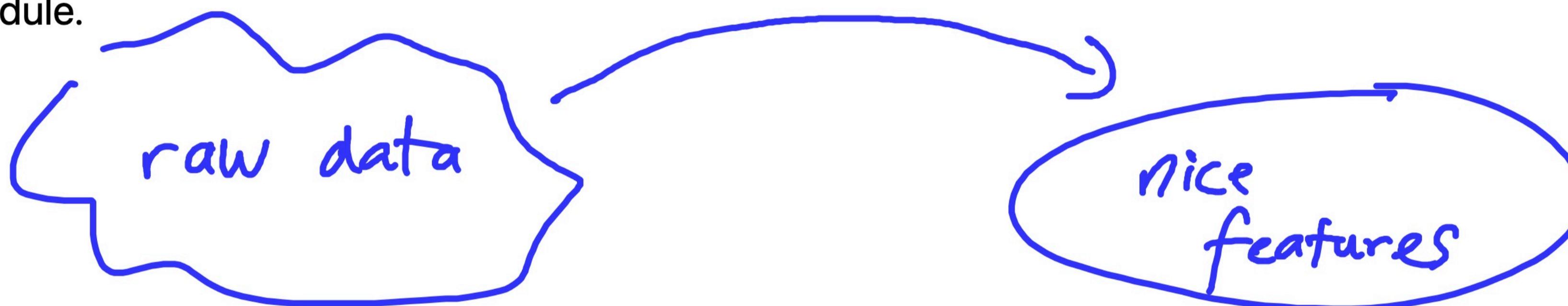
$$T(x_i) = \beta_0 + \beta_1 x_i$$



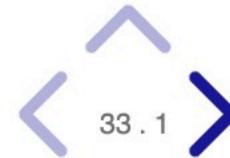
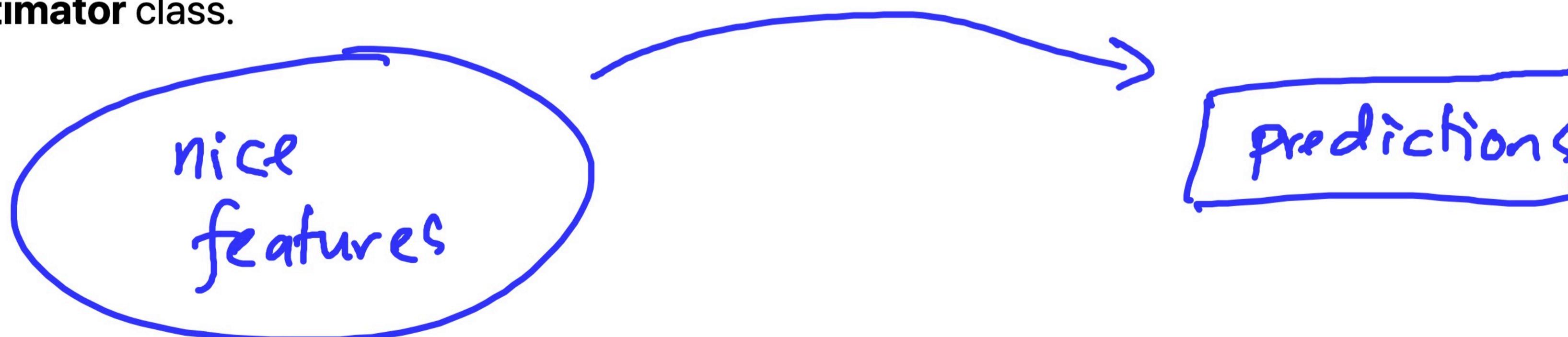


preprocessing and linear_models

- For the **feature engineering** step of the modeling pipeline, we will use `sklearn's preprocessing` module.



- For the **model creation** step of the modeling pipeline, we will use `sklearn's linear_model` module, as we've already seen. `linear_model.LinearRegression` is an example of an **estimator** class.



```
In [44]: 1 sales = pd.read_csv('data/sales.csv')
          2 sales.head()
```

Out [44] :

	net_sales	sq_ft	inventory	advertising	district_size	competing_stores
0	231.0	3.0	294	8.2	8.2	11
1	156.0	2.2	232	6.9	4.1	12
2	10.0	0.5	149	3.0	4.3	15
3	519.0	5.5	600	12.0	16.1	1
4	437.0	4.4	567	10.6	14.1	5

2

- For each of 26 stores, we have:
 - net sales,
 - square feet,
 - inventory,
 - advertising expenditure,
 - district size, and
 - number of competing stores.

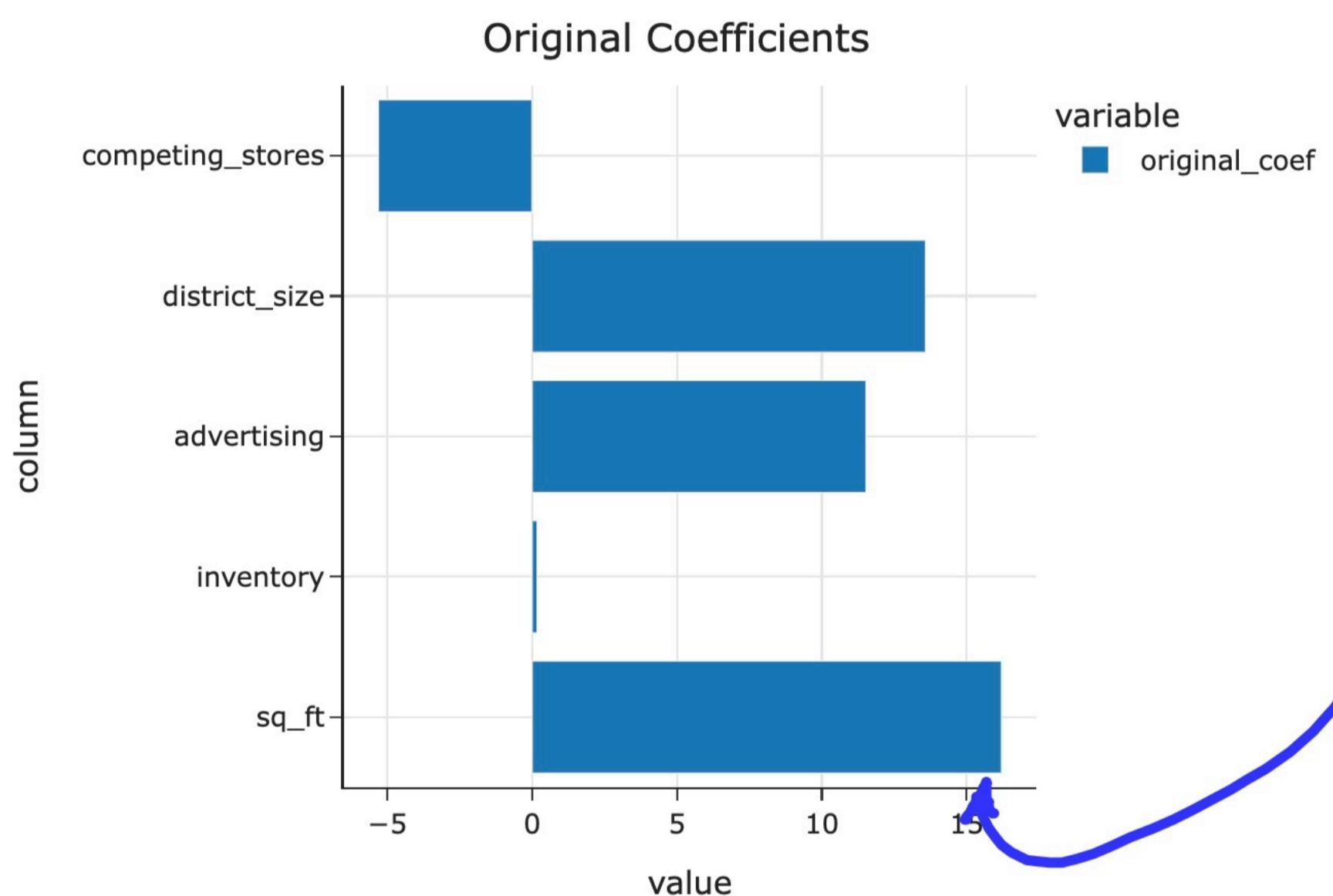
- Our goal is to predict 'net_sales' as a function of other features.



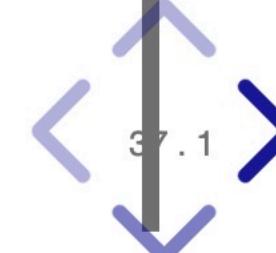
In [47]: 1 sales_model.coef_

Out[47]: array([16.2 , 0.17, 11.53, 13.58, -5.31])

In [48]: 1 coefs = pd.DataFrame().assign(
2 column=sales.columns[1:],
3 original_coef=sales_model.coef_,
4).set_index('column')
5 coefs.plot(kind='barh', title='Original Coefficients')

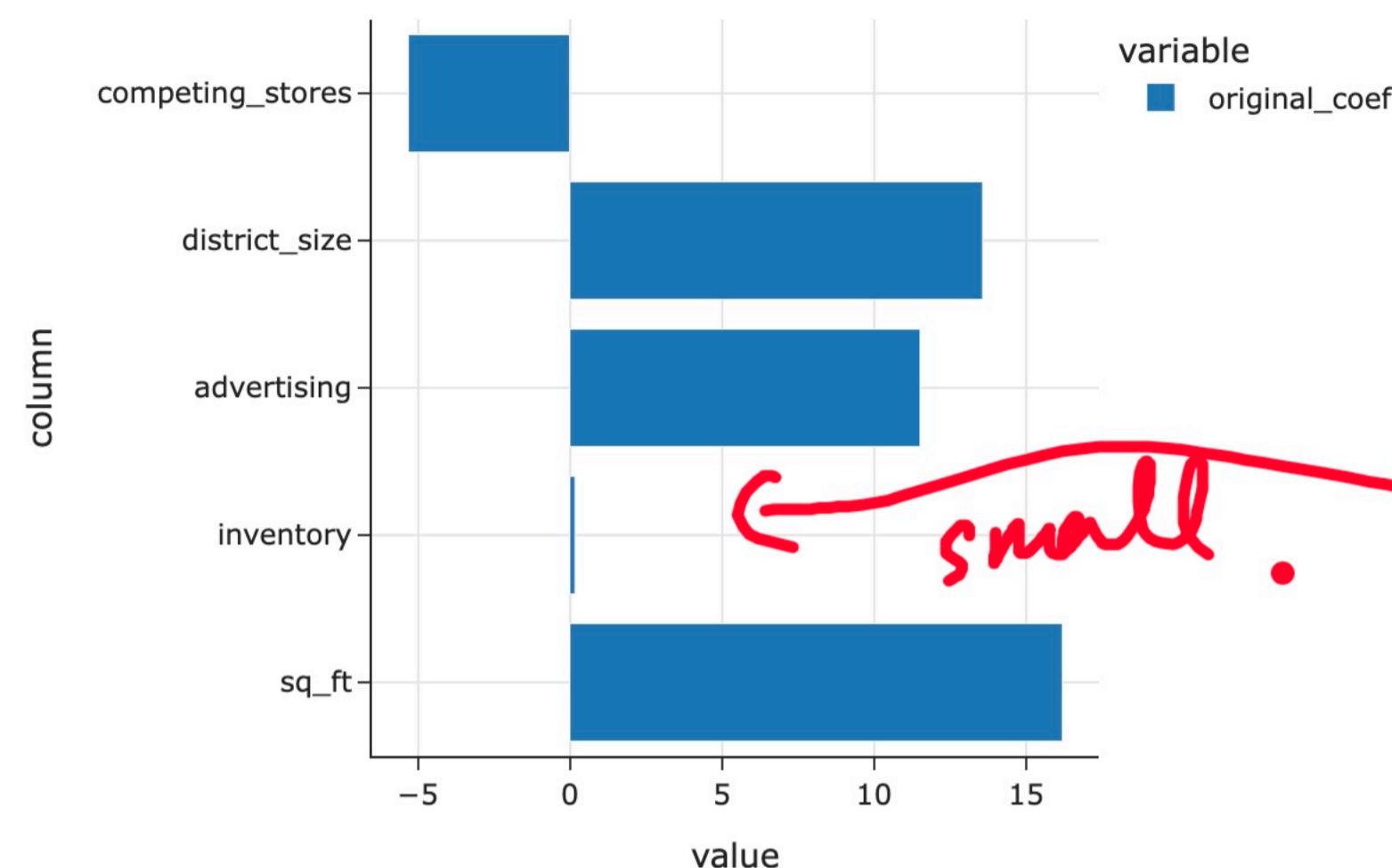


plotting the value for each feature.





Original Coefficients



- What do you notice?

```
In [49]: 1 sales.iloc[:, 1:]
```

```
Out[49]:
```

	sq_ft	inventory	advertising	district_size	competing_stores
0	3.0	294	3.2	8.2	11
1	2.2	232	6.0	4.1	12
2	0.5	149	3.0	4.3	15
...
24	3.5	382	9.8	11.5	5
25	5.1	590	12.0	15.7	0
26	8.6	517	7.0	12.0	8

27 rows × 5 columns

numbers are very big!!!

much smaller scale





features are on
very different scales in.

- What do you notice?

```
In [49]: 1 sales.iloc[:, 1:]
```

```
Out[49]:
```

	sq_ft	inventory	advertising	district_size	competing_stores
0	3.0	294	3.2	8.2	11
1	2.2	232	6.0	4.1	12
2	0.5	149	3.0	4.3	15
...
24	3.5	382	9.8	11.5	5
25	5.1	590	12.0	15.7	0
26	8.6	517	7.0	12.0	8

2.94
2.32
1.49

⋮

much smaller scale

27 rows × 5 columns





Standardization

- When we standardize two or more features, we bring them to the **same scale**.
- Recall: to standardize a feature x_1, x_2, \dots, x_n , we use the formula:

$$z(x_i) = \frac{x_i - \bar{x}}{\sigma_x}$$

do this separately

- Example: 1, 7, 7, 9.

- Mean: $\frac{1+7+7+9}{4} = \frac{24}{4} = 6$.

- Standard deviation:

for each column.

$$\text{SD} = \sqrt{\frac{1}{4}((1-6)^2 + (7-6)^2 + (7-6)^2 + (9-6)^2)} = \sqrt{\frac{1}{4} \cdot 36} = 3$$

- Standardized data:

$$1 \mapsto \frac{1-6}{3} = \boxed{-\frac{5}{3}}$$

$$7 \mapsto \frac{7-6}{3} = \boxed{\frac{1}{3}}$$

$$7 \mapsto \boxed{\frac{1}{3}}$$

$$9 \mapsto \frac{9-6}{3} = \boxed{1}$$



Here, "fitting" the transformer involves computing and saving the mean and SD of each column.

In [53]: 1 stdscaler = StandardScaler()

In [55]: 1 # This is like saying "determine the mean and SD of each column in sales,
2 # other than the 'net_sales' column".
3 stdscaler.fit(sales.iloc[:, 1:])

Out[55]:

StandardScaler
StandardScaler()

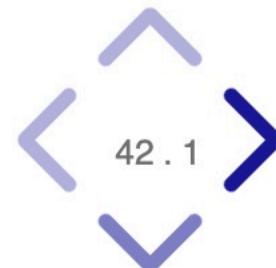
uses mean and SD of
columns it was fit on!

In [57]: 1 stdscaler.transform([[5, 300, 10, 15, 6]])

Out[57]: array([[0.85, -0.47, 0.51, 1.05, -0.36]])

In []: 1 stdscaler.transform(sales.iloc[:, 1:].tail(5))

300 is 0.47 SDs BELow
the average inventory in
Sales:



```
In [55]: 1 stdscaler = StandardScaler()
```

```
In [55]: 1 # This is like saying "determine the mean and SD of each column in sales,  
2 # other than the 'net_sales' column".  
3 stdscaler.fit(sales.iloc[:, 1:])
```

Out[55]:

▼ StandardScaler ⓘ ⓘ
StandardScaler()

```
In [57]: 1 stdscaler.transform([[5, 300, 10, 15, 6]])
```

Out[57]: array([[0.85, -0.47, 0.51, 1.05, -0.36]])

```
In [59]: 1 stdscaler.transform(sales.iloc[:, 1:].tail(5))
```

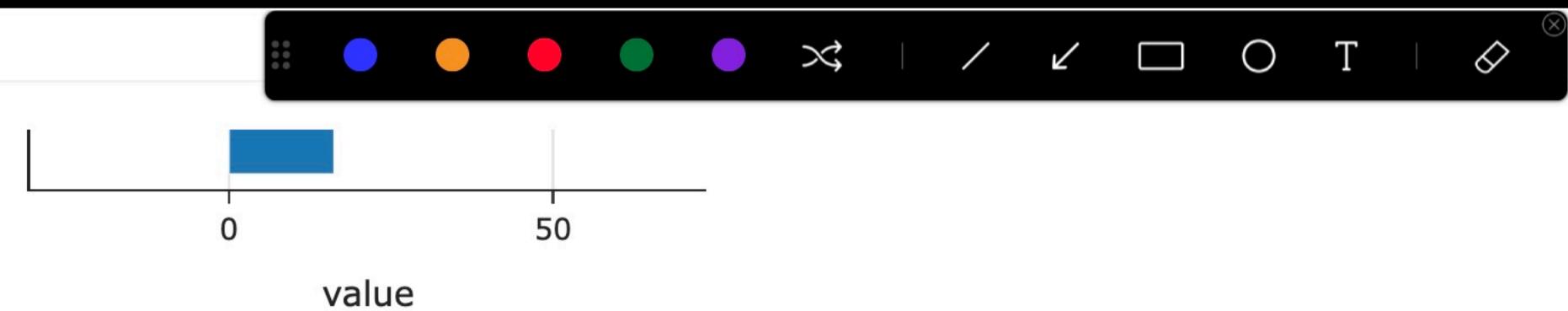
Out[59]: array([[-1.13, -1.31, -1.35, -1.6, 0.89],
[0.14, 0.39, 0.4, 0.32, -0.36],
[0.09, -0.03, 0.46, 0.36, -0.57],
[0.9, 1.08, 1.05, 1.19, -1.61],
[2.67, 0.69, -0.3, 0.46, 0.05]])

```
In [60]: 1 new_scaler = StandardScaler()  
2 new_scaler.fit(sales.iloc[:, 1:].tail(5))  
3 new_scaler.transform(sales.iloc[:, 1:].tail(5))
```

Out[60]: array([[-1.33, -1.79, -1.71, -1.88, 1.48],
[-0.32, 0.28, 0.43, 0.19, -0.05],
[-0.36, -0.24, 0.49, 0.23, -0.31],
[0.29, 1.11, 1.22, 1.13, -1.58],
[1.71, 0.64, -0.43, 0.34, 0.46]])

Why different?
because different
datasets were
used in fit!





- Did the performance of the resulting model change?

```
In [67]: 1 mean_squared_error(sales.iloc[:, 0],  
2 sales_model.predict(sales.iloc[:, 1:]))
```

Out[67]: 242.27445717154964

```
In [69]: 1 sales_model.predict(sales.iloc[:, 1:])
```

Out[69]: array([228.54, 128.78, 28.57, ..., 347.13, 518.33, 411.92])

same preds!

```
In [68]: 1 mean_squared_error(sales.iloc[:, 0],  
2 sales_model_std.predict(stdscaler.transform(sales.iloc[:, 1:])))
```

Out[68]: 242.27445717154956

```
In [70]: 1 sales_model_std.predict(stdscaler.transform(sales.iloc[:, 1:]))
```

Out[70]: array([228.54, 128.78, 28.57, ..., 347.13, 518.33, 411.92])

- No!

The span of the design matrix did not change, so the predictions did not change. It's just the coefficients that changed.

$$X = \begin{bmatrix} \vec{x}^{(1)} \\ \vdots \\ \vec{x}^{(n)} \end{bmatrix}$$

$$\begin{aligned}
 \vec{z}^{(1)} &= \frac{\vec{x}^{(1)} - \text{mean } (\vec{x}^{(1)})}{SD(\vec{x}^{(1)})} \\
 &= \frac{\vec{x}^{(1)}}{SD(\vec{x}^{(1)})} - \frac{\text{mean } (\vec{x}^{(1)})}{SD(\vec{x}^{(1)})} \\
 &= \frac{1}{SD(\vec{x}^{(1)})_a} \vec{x}^{(1)} - \frac{\text{mean } (\vec{x}^{(1)})}{SD(\vec{x}^{(1)})_b}
 \end{aligned}$$

A hand-drawn diagram featuring large blue cursive text "big picture" at the top left. A thick blue curved line starts from the bottom left, goes up and to the right, then turns downwards towards the bottom right. Below the curve, the word "big" is written in orange cursive, with a large orange "t" partially obscured by the blue line. To the right of "big", the word "is" is written in orange. At the bottom right, there is orange cursive text "in. contr. of" followed by a large orange "X".

$$\vec{v}_1 = a \vec{x}^{(1)} + b \vec{1}$$