

Lecture 9

Multiple Linear Regression

DSC 40A, UCSD

Agenda

- Multiple linear regression.
- Interpreting parameters.
- Feature engineering and transformations.

Question 🤔

Take a moment to pause and reflect...

If you have any questions please post online to our forms/Q&A site.

Course staff will answer them ASAP!

Multiple linear regression

features

	departure_hour	day_of_month	minutes
0	10.816667	15	68.0
1	7.750000	16	94.0
2	8.450000	22	63.0
3	7.133333	23	100.0
4	9.150000	30	69.0
...

So far, we've fit **simple** linear regression models, which use only **one** feature
 ('departure_hour') for making predictions.

Incorporating multiple features

- In the context of the commute times dataset, the simple linear regression model we fit was of the form:

$$\begin{aligned}\text{pred. commute} &= H(\text{departure hour}) \\ &= w_0 + \underline{w_1 \cdot \text{departure hour}} \quad \text{input feature}\end{aligned}$$

- Now, we'll try and fit a multiple linear regression model of the form:

$$\begin{aligned}\text{pred. commute} &= H(\text{departure hour}) \\ &= w_0 + \underline{w_1 \cdot \text{departure hour}} + \underline{w_2 \cdot \text{day of month}} \quad \text{two input features}\end{aligned}$$

- Linear regression with **multiple** features is called **multiple linear regression**.

- How do we find w_0^* , w_1^* , and w_2^* ? \Rightarrow via normal equations

Geometric interpretation

- The hypothesis function:

$$H(\text{departure hour}) = w_0 + \underline{w_1} \cdot \text{departure hour}$$

looks like a **line** in 2D.

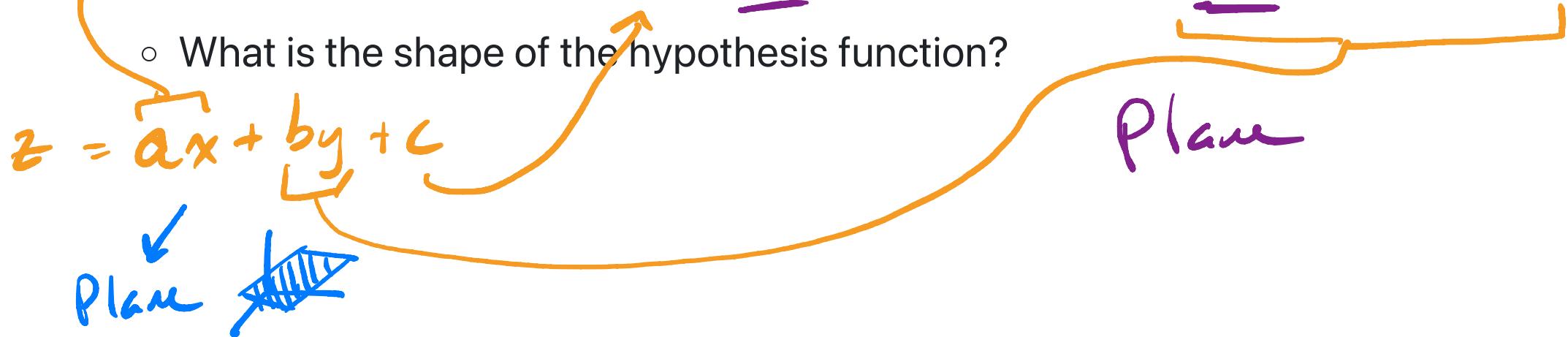
line

- Questions:

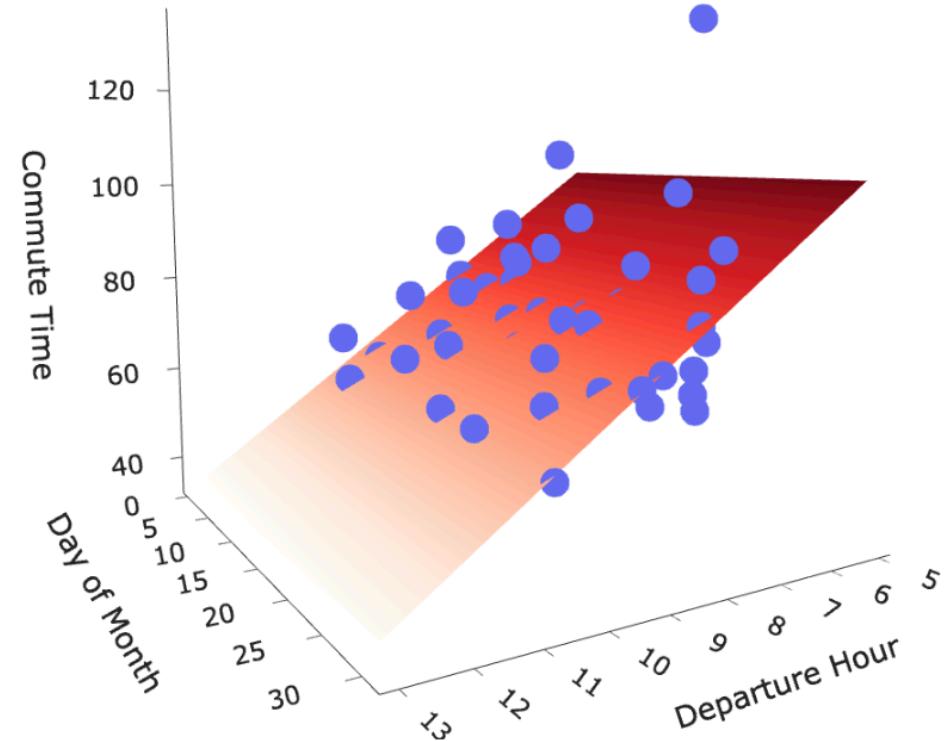
- How many dimensions do we need to graph the hypothesis function:

$$H(\text{departure hour}) = w_0 + \underline{w_1} \cdot \text{departure hour} + \underline{w_2} \cdot \text{day of month}$$

- What is the shape of the hypothesis function?



Commute Time vs. Departure Hour and Day of Month



Our new hypothesis function is a **plane** in 3D!

Our goal is to find the **plane** of best fit that pierces through the cloud of points.

The setup

- Suppose we have the following dataset.

Handwritten annotations:

- xs* (blue wavy line) points to the first two columns: `departure_hour` and `day_of_month`.
- ys* (orange wavy line) points to the last column: `minutes`.
- row* (black wavy line) points to the vertical axis on the left.
- Observations* (black handwritten text) is written vertically next to the row index.
- "Predictions"* (orange handwritten text) is written to the right of the `minutes` column.

row	Observations	departure_hour	day_of_month	minutes	"Predictions"
1		8.45	22	63.0	
2		8.90	28	89.0	
3		8.72	18	89.0	

- We can represent each day with a feature vector, \vec{x} :

$$\vec{x}_1 = \begin{bmatrix} 8.45 \\ 22 \end{bmatrix}$$

$$\vec{x}_2 = \begin{bmatrix} 8.90 \\ 28 \end{bmatrix}$$

$$\vec{x}_3 = \begin{bmatrix} 8.72 \\ 18 \end{bmatrix}$$

The hypothesis vector

- When our hypothesis function is of the form:

$$H(\text{departure hour}) = w_0 + w_1 \cdot \text{departure hour} + w_2 \cdot \text{day of month}$$

the hypothesis vector $\vec{h} \in \mathbb{R}^n$ can be written as:

$$\vec{h} = \begin{bmatrix} H(\text{departure hour}_1, \text{day}_1) \\ H(\text{departure hour}_2, \text{day}_2) \\ \dots \\ H(\text{departure hour}_n, \text{day}_n) \end{bmatrix} = \begin{bmatrix} 1 & \text{departure hour}_1 & \text{day}_1 \\ 1 & \text{departure hour}_2 & \text{day}_2 \\ \dots & \dots & \dots \\ 1 & \text{departure hour}_n & \text{day}_n \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

$n \times 3 \cdot 3 \times 1$

feature Vector \vec{x}_i

$$= \begin{bmatrix} \text{dep. hour}_i \\ \text{day}_i \end{bmatrix}$$

\vec{x}_i^T

$$\vec{h} = \vec{X} \vec{\omega}$$

All Predictions

Design Matrix

Parameter Vector

Finding the optimal parameters

- To find the optimal parameter vector, \vec{w}^* , we can use the **design matrix** $X \in \mathbb{R}^{n \times 3}$ and **observation vector** $\vec{y} \in \mathbb{R}^n$:

$$X = \begin{bmatrix} 1 & \text{departure hour}_1 & \text{day}_1 \\ 1 & \text{departure hour}_2 & \text{day}_2 \\ \dots & \dots & \dots \\ 1 & \text{departure hour}_n & \text{day}_n \end{bmatrix} \quad \vec{y} = \begin{bmatrix} \text{commute time}_1 \\ \text{commute time}_2 \\ \vdots \\ \text{commute time}_n \end{bmatrix}$$

- Then, all we need to do is solve the **normal equations**:

$$X^T X \vec{w}^* = X^T \vec{y}$$

If $X^T X$ is invertible, we know the solution is:

$$\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

Notation for multiple linear regression

- We will need to keep track of multiple features for every individual in our dataset.
 - In practice, we could have hundreds or thousands of features!
- As before, subscripts distinguish between individuals in our dataset. We have n individuals, also called **training examples**.
- Superscripts distinguish between **features**. We have d features.

$x^{(1)}, x^{(2)}, \dots, x^{(d)}$

Think of $x^{(1)}, x^{(2)}, \dots$ as new variable names, like new letters.

departure hour: $x^{(1)}$

day of month: $x^{(2)}$

x^2 ← exponent
 $x^{(2)}$ ← not an exponent

$x_4^{(5)}$

represents the value in our dataset
if the 5th feature & fourth training example

Augmented feature vectors

- The **augmented feature vector** $\text{Aug}(\vec{x})$ is the vector obtained by adding a 1 to the front of feature vector \vec{x} :

$x^{(1)} = \text{Dep. Hour}$
 $x^{(2)} = \text{day of Month}$
 $x^{(3)} = \text{Avg tire pres.}$
 $x^{(4)} = \text{time spent eating Breakfast}$

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

$$\text{Aug}(\vec{x}) = \begin{bmatrix} 1 \\ x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{bmatrix}_{(d+1) \times 1}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}_{(d+1) \times 1}$$

- Then, our hypothesis function is:

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

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:

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

- We want to find a good linear hypothesis function:

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

$\{w_0^*, w_1^*, w_d^*\}$ find these

The general solution

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

a single training example + 1 = all the "stuff" we need to make a prediction

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

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

n × (d+1)

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

$$X^T X \vec{w}^* = X^T \vec{y}$$

a single feature

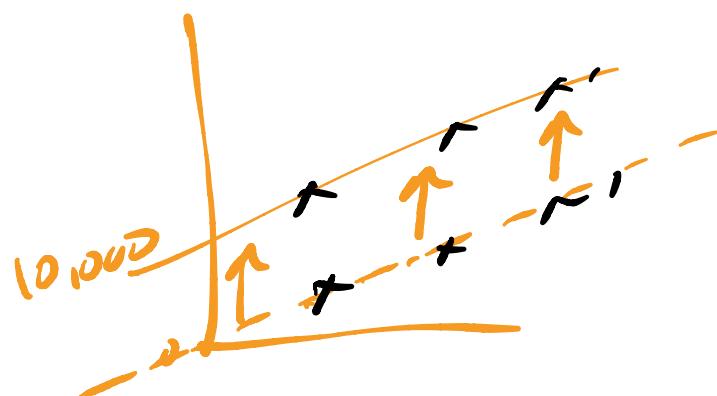
e.g. "departure hour"

Terminology for parameters

- With d features, \vec{w} has $d + 1$ entries.
- w_0 is the bias, also known as the intercept.
- w_1, w_2, \dots, w_d each give the weight, or coefficient, or slope, of a feature.

$$H(\vec{x}) = w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)}$$

if all weights were 0
bias would be our pred.



Interpreting parameters

Example: Predicting sales

- For each of 26 stores, we have:

y {

- net sales,
- square feet,
- inventory,
- advertising expenditure,
- district size, and
- number of competing stores.

X_s {

$$n = 26$$

$$d = 5$$

not a feature

$$\vec{w} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 10 \\ 2.5 \\ 0.21 \end{bmatrix}$$

- Goal:** Predict net sales given the other five features.
- To begin, we'll start trying to fit the hypothesis function to predict sales:

$$H(\text{square feet, competitors}) = \underline{w_0 + w_1 \cdot \text{square feet}} + \underline{w_2 \cdot \text{competitors}}$$

$d = 2$

Which features are most "important"?

- The most important feature is **not necessarily** the feature with largest magnitude weight.
- Features are measured in different units, i.e. different scales.
 - Suppose I fit one hypothesis function, H_1 , with sales in US dollars, and another hypothesis function, H_2 , with sales in Japanese yen ($1 \text{ USD} \approx 157 \text{ yen}$).
 - Sales is just as important in both hypothesis functions.
 - But the weight of sales in H_1 will be 157 times larger than the weight of sales in H_2 .
- **Solution:** If you care about the interpretability of the resulting weights, **standardize** each feature before performing regression, i.e. convert each feature to standard units.

Standard units

- Recall: to convert a feature x_1, x_2, \dots, x_n to standard units, we use the formula:

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

- Example: 1, 7, 7, 9.

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

- Standard deviation:

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

- Standardized data:

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

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

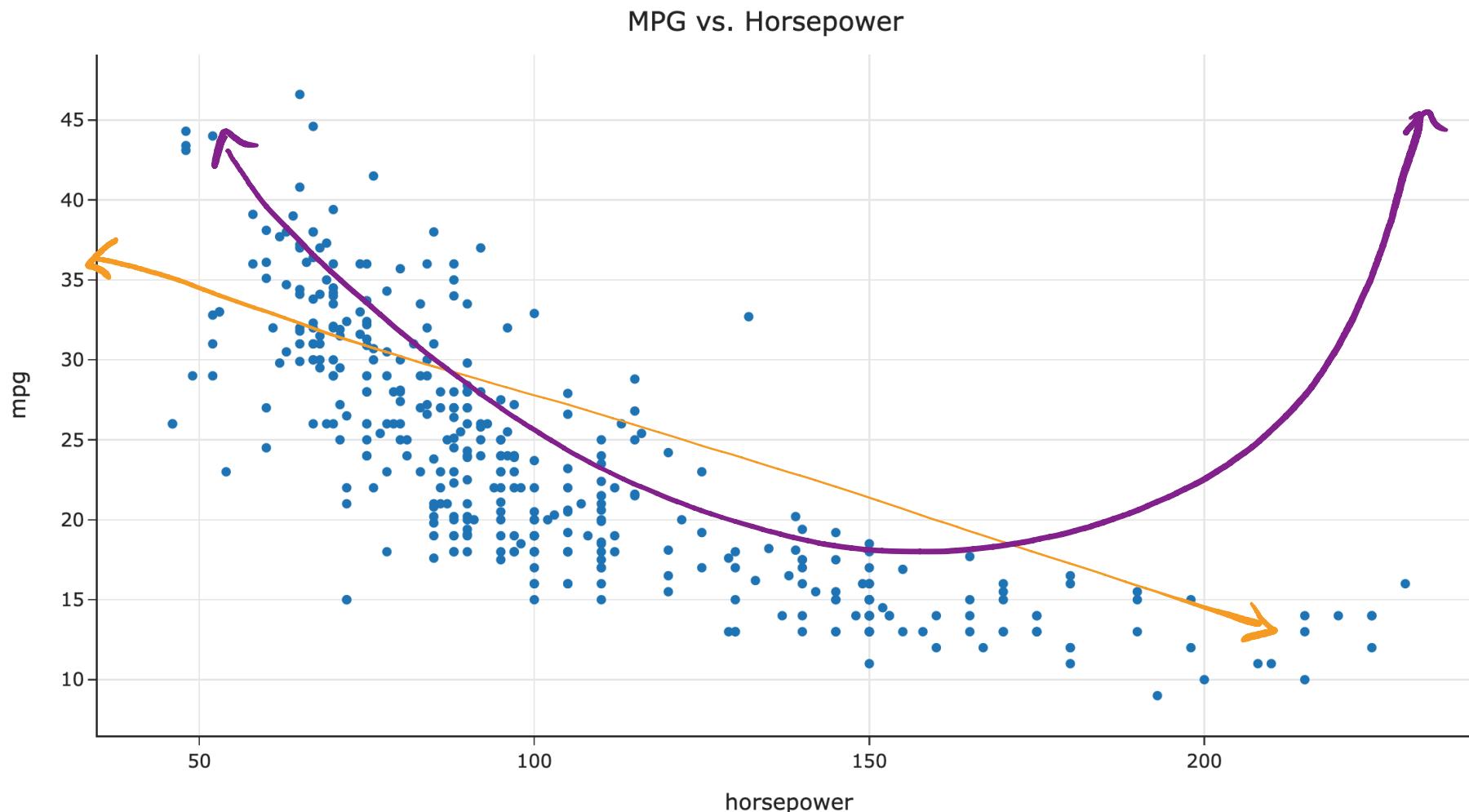
$$7 \mapsto \boxed{\frac{1}{3}}$$

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

Standard units for multiple linear regression

- The result of standardizing each feature (separately!) is that the units of each feature are on the same scale.
 - There's no need to standardize the outcome (net sales), since it's not being compared to anything.
 - Also, we can't standardize the column of all 1s.
- Then, solve the normal equations. The resulting $w_0^*, w_1^*, \dots, w_d^*$ are called the **standardized regression coefficients**.
- Standardized regression coefficients can be directly compared to one another.
- Note that standardizing each feature **does not** change the MSE of the resulting hypothesis function!

Feature engineering and transformations



Question: Would a linear hypothesis function work well on this dataset?

$$\omega_2 x^2 \neq \omega_2 x^{(2)}$$

A quadratic hypothesis function

- It looks like there's some sort of quadratic relationship between horsepower and MPG in the last scatter plot. We want to try and fit a hypothesis function of the form:

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

$$H(x) = \omega_0 + \omega_1 \cdot \underline{\quad} + \omega_2 \cdot \underline{\quad}$$

↑
no ω_3

- Note that while this is quadratic in horsepower, it is **linear in the parameters!**
- That is, it is a **linear combination of features**.
- We can do that, by choosing our two "features" to be x_i and x_i^2 , respectively.
 - In other words, $x_i^{(1)} = x_i$ and $x_i^{(2)} = x_i^2$.
 - More generally, we can create new features out of existing features.

Feature
Engineering

$$x_i^{(2)} = x_i^2 = HP \text{ (squared)}$$

A quadratic hypothesis function

- Desired hypothesis function: $H(x) = w_0 + w_1x + w_2x^2$.
- The resulting design matrix looks like:

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \dots & & \\ 1 & x_n & x_n^2 \end{bmatrix} = \begin{bmatrix} 1 & h_{p,1} & h_{p,1}^2 \\ 1 & h_{p,2} & h_{p,2}^2 \\ \vdots & \vdots & \vdots \\ 1 & h_{p,n} & h_{p,n}^2 \end{bmatrix}_{n \times 3}$$

- To find the optimal parameter vector \vec{w}^* , we need to solve the **normal equations**!

$$X^T X \vec{w}^* = X^T \vec{y}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

More examples

- What if we want to use a hypothesis function of the form:
 $H(x) = w_0 + w_1x + w_2x^2 + w_3x^3?$

Design Matrix

$(n \times 4)$ $\vec{w} (4 \times 1)$



This is fine

- What if we want to use a hypothesis function of the form:

$$H(x) = w_1 \frac{1}{x^2} + w_2 \sin x + w_3 e^x?$$

1 2 3

w_0

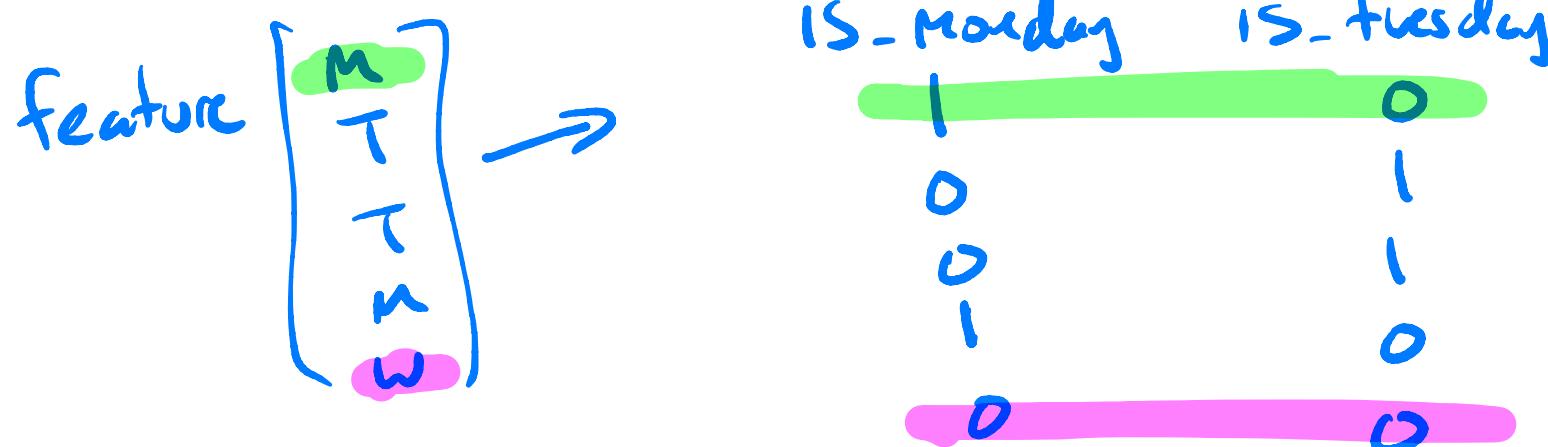
$$X = \begin{bmatrix} \frac{1}{x_1^2} & \sin x_1 & e^{x_1} \\ \frac{1}{x_2^2} & \sin x_2 & e^{x_2} \\ \vdots & \vdots & \vdots \\ \frac{1}{x_n^2} & \sin x_n & e^{x_n} \end{bmatrix}_{(n \times 3)}$$

iff $X \vec{w} = \vec{y}$

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

Feature engineering

- The process of creating new features out of existing information in our dataset is called **feature engineering**.
- In this class, feature engineering will mostly be restricted to creating non-linear functions of existing features (as in the previous example).
- In the future you'll learn how to do other things, like encode categorical information.
 - You'll be exposed to this in Homework 4, Problem 5!



Non-linear functions of multiple features

- Recall our earlier example of predicting sales from square footage and number of competitors. What if we want a hypothesis function of the form:

$$\begin{aligned}H(\text{sqft}, \text{comp}) &= w_0 + w_1 \cdot \text{sqft} + w_2 \cdot \text{sqft}^2 + w_3 \cdot \text{comp} + w_4 \cdot (\text{sqft} \cdot \text{comp}) \\&= w_0 + w_1 s + w_2 s^2 + w_3 c + w_4 sc\end{aligned}$$

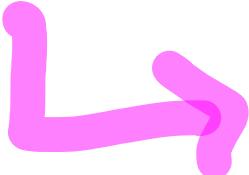
- The solution is to choose a design matrix accordingly:

$$X = \begin{bmatrix} 1 & s_1 & s_1^2 & c_1 & s_1 c_1 \\ 1 & s_2 & s_2^2 & c_2 & s_2 c_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & s_n & s_n^2 & c_n & s_n c_n \end{bmatrix}$$

$$\frac{w_0 + w_1 - + w_2 - + \cdots + w_n -}{L.C. \quad L.C. \quad L.C. \quad L.C.}$$

Finding the optimal parameter vector, \vec{w}^*

- As long as the form of the hypothesis function permits us to write $\vec{h} = \vec{X}\vec{w}$ for some \vec{X} and \vec{w} , the mean squared error is:

 $R_{\text{sq}}(\vec{w}) = \frac{1}{n} \|\vec{y} - \vec{X}\vec{w}\|^2$

Design matrix Param. vector

- Regardless of the values of X and \vec{y} , the value of \vec{w}^* that minimizes $R_{\text{sq}}(\vec{w})$ is the solution to the normal equations:

$$X^T X \vec{w}^* = X^T \vec{y}$$

if invertible
 λ
 $X^T X$ is then Sol.
 for \vec{w}^*

Linear in the parameters

- We can fit rules like:

$$w_0 + w_1 x + w_2 x^2$$

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

- This includes arbitrary polynomials.
- These are all linear combinations of (just) features.

- We can't fit rules like:

X $w_0 + e^{w_1 x}$

w is an exponent
the w_s are in this sin()

X $w_0 + \sin(w_1 x^{(1)} + w_2 x^{(2)})$

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

these w_s are Dirty

- We can have any number of parameters, as long as our hypothesis function is **linear in the parameters**, or linear when we think of it as a function of the parameters.