

EECS 245, Spring 2026

LEC 9

Regression using Linear Algebra;
The Gradient Vector

→ Read: Ch. 7 - Ch. 8

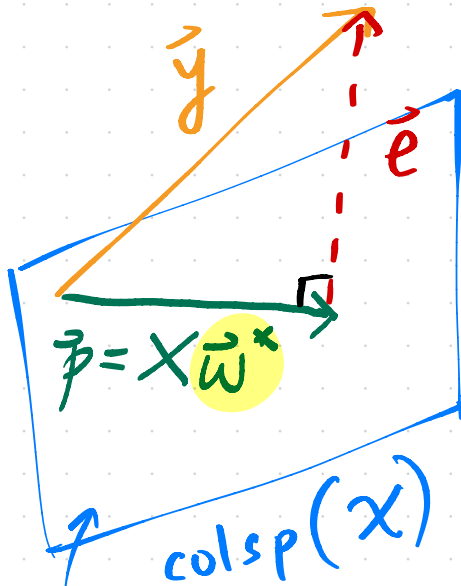
Agenda

- Idea 1: How does any of this relate to linear regression and ML?
 - Projections
 - The modeling recipe
 - The "design matrix"
 - Using multiple features
- Idea 2: Gradients
 - Calculus meets linear algebra

Announcements

- Due dates:
 - Lab 8: tomorrow
 - HW 7: Thursday
 - HW 8: Sunday. no slip days!
- Lab 7, HW 6 solutions posted
- MT 2 on Tuesday, June 9th, 1-3PM: see logistics at eecs245.org/mt2

Recap: Projections



every vector in $\text{colsp}(X)$
is of form $X\vec{v}$, $\vec{v} \in \mathbb{R}^d$

X : $n \times d$ matrix
(d columns, each in \mathbb{R}^n)
 $\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(d)}$

$\vec{y} \in \mathbb{R}^n$

Q: What vector in $\text{colsp}(X)$
is closest to \vec{y} ?

A: The vector $X\vec{w}^*$, where
 $\vec{y} - X\vec{w}^*$ is orthogonal
error to $\text{colsp}(X)$.

How did we find \vec{w}^* ?

It is the solution to the normal equation

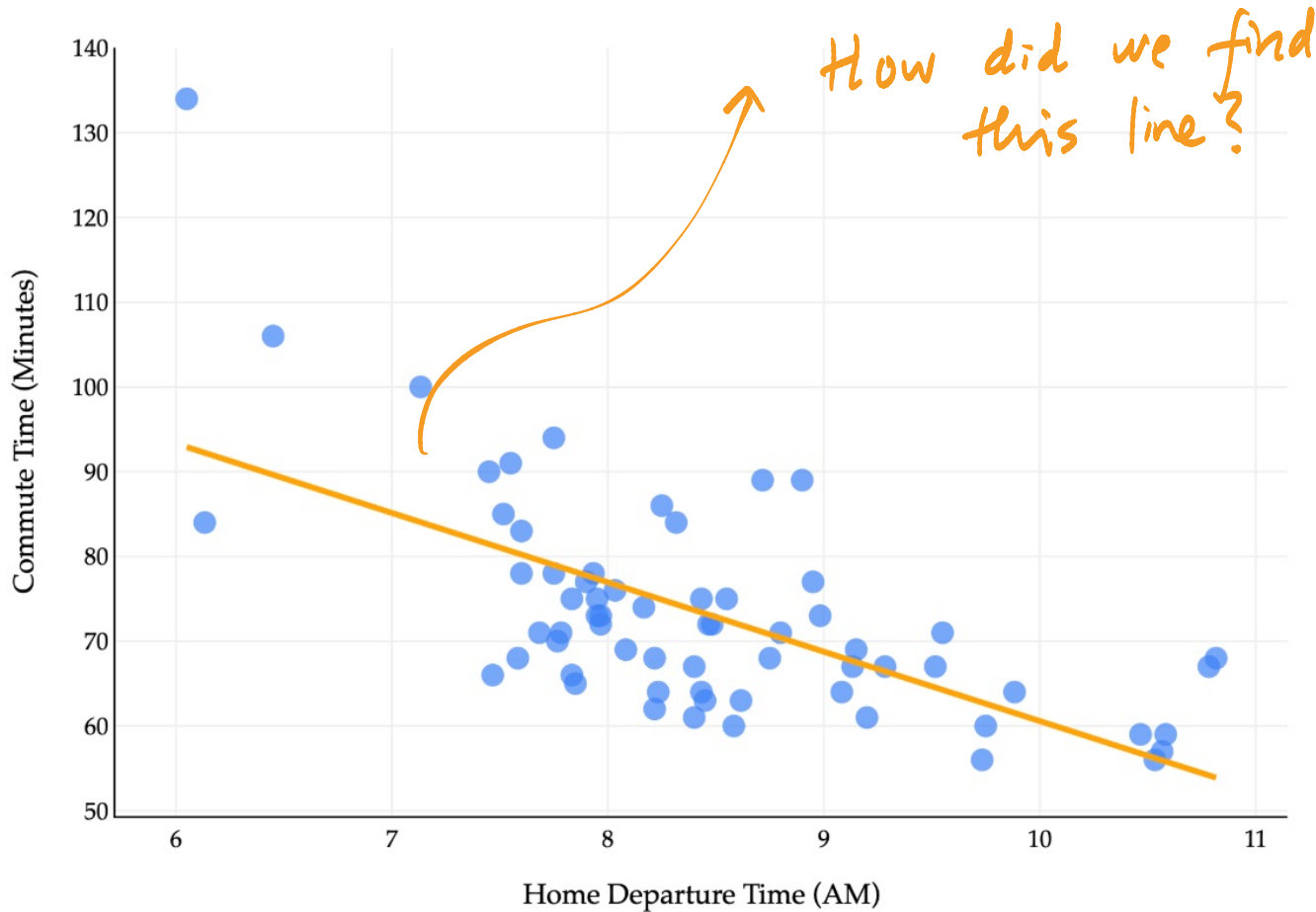
$$X^T X \vec{w} = X^T \vec{y}$$

① If X 's cols are linearly independent, then

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

unique!

② otherwise, there are infinitely many best \vec{w}^* 's!
→ see HW 7, Problem 3



Recap: Three-step modeling recipe (empirical risk minimization)

① Choose a model

$$h(x_i) = w_0 + w_1 x_i$$

② Choose a loss function

$$L_{sq}(y_i, h(x_i)) = \underbrace{(y_i - h(x_i))^2}_{(\text{actual-pred})^2}$$

③ Minimize average loss to find optimal parameters

mean squared error

$$R_{sq}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n \left(y_i - (w_0 + w_1 x_i) \right)^2$$

looks like error!

→ before: to find w^* 's, we took partial derivatives

→ if we had more w 's (more features), this might get too complicated

→ maybe there's a way to rewrite R_{sq} using linear algebra!

→ goal: "convert" mean squared error

$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$ actual commute times
"observation vector"

$\vec{\hat{y}} \in \mathbb{R}^n$

predicted commute times
"prediction vector"

$$\vec{\hat{p}} = \begin{bmatrix} h(x_1) \\ h(x_2) \\ \vdots \\ h(x_n) \end{bmatrix} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ \vdots \\ w_0 + w_1 x_n \end{bmatrix}$$

"design matrix"

X

"parameter vector"

\vec{w}

$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{matrix} \\ \\ \\ \\ \\ \end{matrix} \begin{matrix} \\ \\ \\ \\ \\ \end{matrix}$$

$n \times 2$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix}_{2 \times 1}$$

before: $R_{sq}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$

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

These formulas
are equivalent!

error vector

$$\begin{bmatrix} y_1 - (w_0 + w_1 x_1) \\ y_2 - (w_0 + w_1 x_2) \\ \vdots \\ y_n - (w_0 + w_1 x_n) \end{bmatrix}$$

We need to find the \vec{w} that minimizes

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

How? We already solved this problem via projections!

The "best" \vec{w} is the \vec{w} such that

$X\vec{w}^*$ is the closest vector in $\text{colsp}(X)$ to \vec{y}

→ If X 's cols are linearly independent, then

$$\vec{w}^* = \begin{bmatrix} w_0^* \\ w_1^* \end{bmatrix} = (X^T X)^{-1} X^T \vec{y}$$

$$\vec{w}^* = (X^T X)^{-1} X^T \vec{y} = \begin{bmatrix} \bar{y} - w_1^* \bar{x} \\ r \frac{\sigma_y}{\sigma_x} \end{bmatrix}$$

These formulas are the same!

Q: Why is it guaranteed that the model's errors add up to 0?

A: Remember, when projecting \hat{y} onto $\text{colsp}(X)$, the error vector \vec{e} is orthogonal to $\text{colsp}(X)$.

In regression, the design matrix X has a column of all 1's. So, it must be that $\vec{e} \cdot \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} = 0 \Rightarrow e_1 + e_2 + \dots + e_n = 0!$

The point was to unlock the ability to add more features!

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

15j

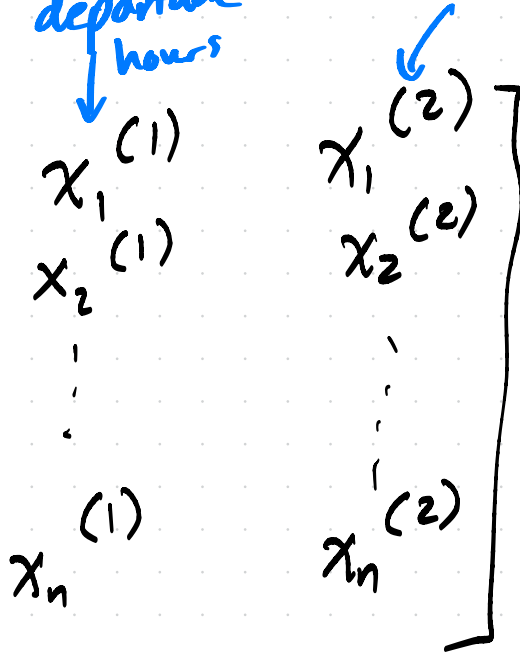
Let's start by adding just one new feature. *not an exponent!*

$$h\left(\underbrace{x_i^{(1)}}_{\text{departure hour}}, \underbrace{x_i^{(2)}}_{\text{day of month}}\right)$$

$$X = \begin{matrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{matrix} \quad n \times 3$$

$$= w_0 + w_1 x_i^{(1)} + w_2 x_i^{(2)}$$

departure hours *days of month*



key idea:
new feature
⇒ new col
in design
matrix!

How to find w_0^* , w_1^* , w_2^* ?

Just use

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

Aside: what if some features are
linear combinations of other features?

⇒ Then, $(X^T X)^{-1}$ doesn't exist

⇒ "Multicollinearity"

⇒ Doesn't impact model performance, but interpretability
and uniqueness of \vec{w}^*

A more general notation for multiple linear regressions

In general, suppose we have a dataset with
 n rows, d features \vec{x}_i "feature vector"

Model:

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

$$= \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \cdot \begin{bmatrix} 1 \\ x_i^{(1)} \\ \vdots \\ x_i^{(d)} \end{bmatrix} = \vec{w} \cdot \text{Aug}(\vec{x}_i)$$

$$\vec{w} \in \mathbb{R}^{d+1}$$

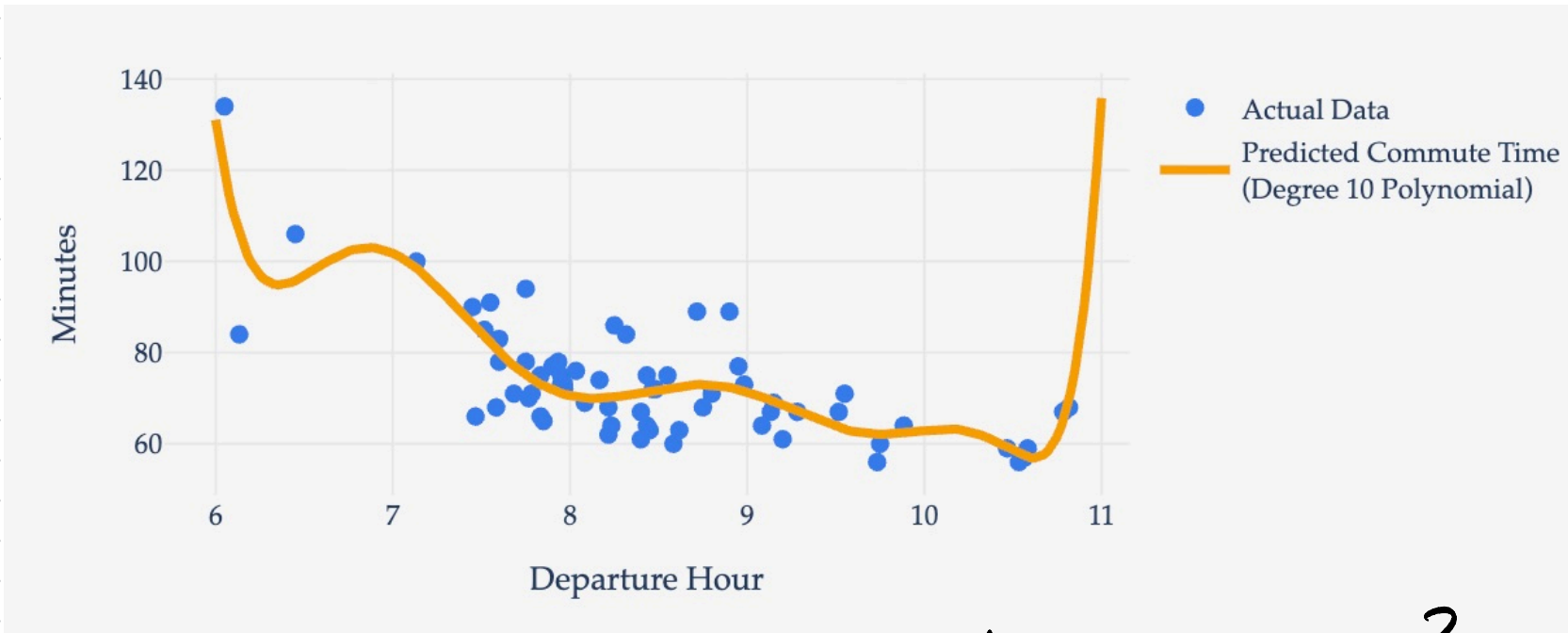
$$\vec{x}_i \in \mathbb{R}^d$$

$$\text{Aug}(\vec{x}_i) \in \mathbb{R}^{d+1}$$

General design matrix has n rows and $d+1$ columns

$$X = \begin{bmatrix} | & x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(d)} \\ | & x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(d)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ | & x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(d)} \end{bmatrix}$$

Find optimal parameters, \vec{w}^* , using normal equations
(which is the same as minimizing mean squared error!)



Can we make this using multiple linear regression?
Sure can! Just make an $n \times 11$ design matrix.

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

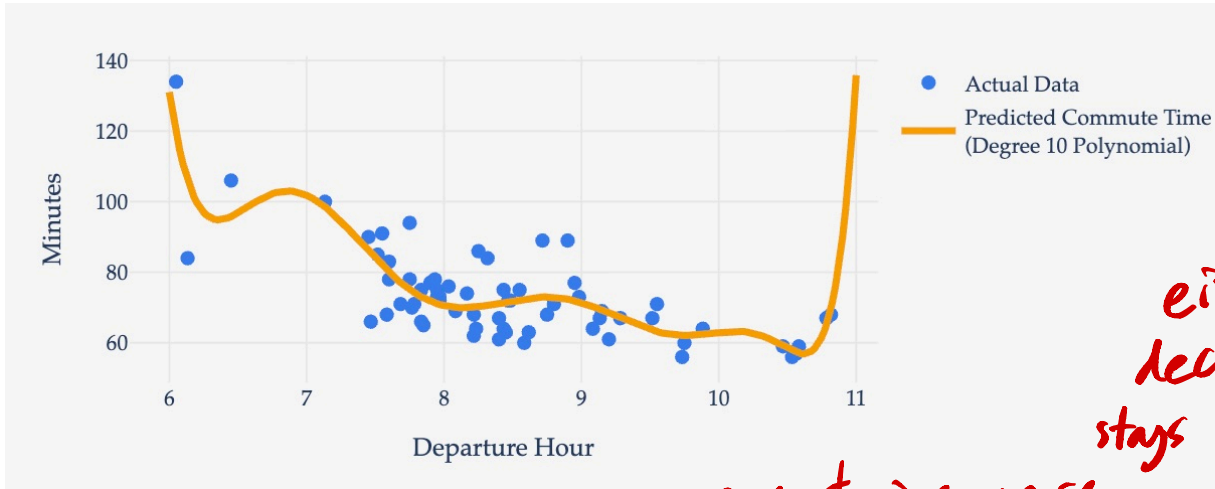
"feature engineering"
 → new features
 from old

degree 3 poly,
 $x_i = \text{dept hour}$

Goal: $h(x_i) = w_0 + w_1 x_i + w_2 x_i^2 + w_3 x_i^3$

How? Write the first 2 rows of design matrix.
 No variables!

$$X = \begin{bmatrix} 1 & 10.81 & 10.81^2 & 10.81^3 \\ 1 & 7.75 & 7.75^2 & 7.75^3 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$



either
decreases or
stays same

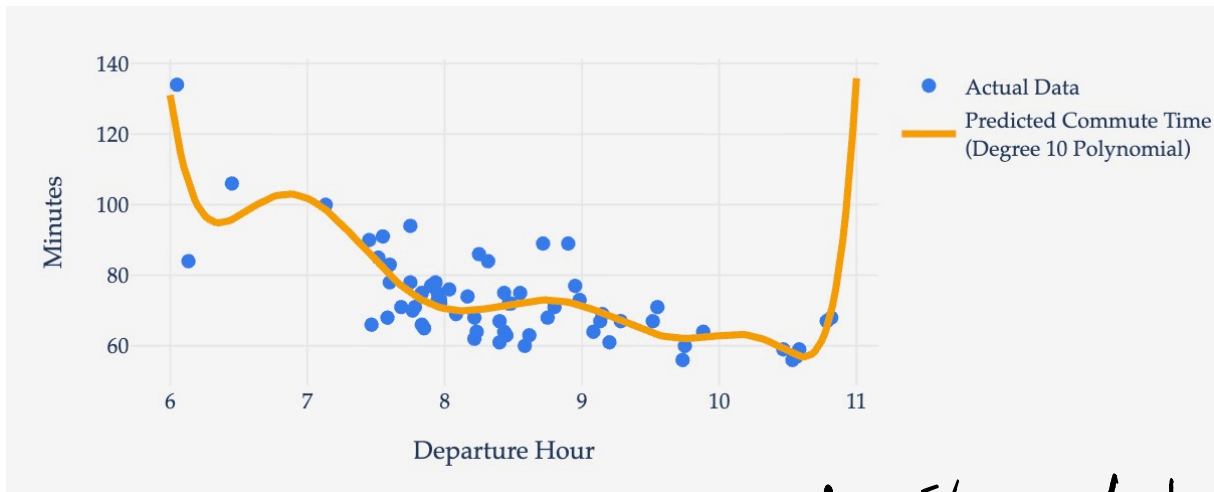
→ Increasing polynomial degree
mean squared error

cannot increase
on the training data

→ But, we care about

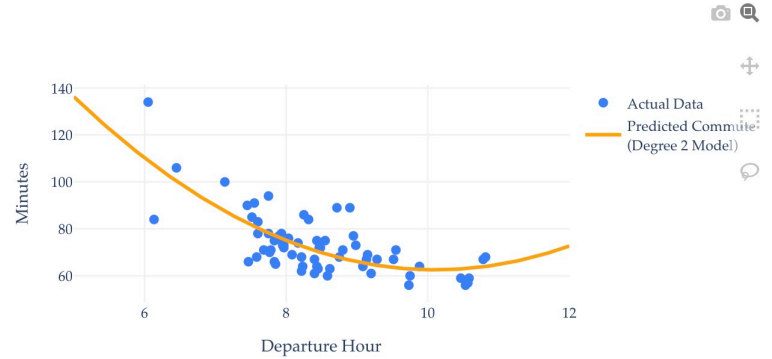
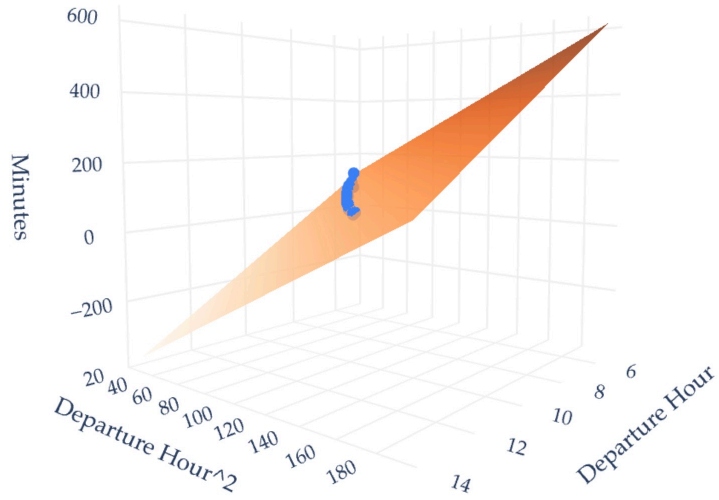
models generalizing to unseen data

→ We need a technique for selecting model complexity
→ see HW 9/10



How is this still linear regression, despite not being a line?

⇒ If we drew the model in \mathbb{R}^n , it would look like a flat surface ("affine" subspace - not guaranteed to go through origin)



What sorts of models can we not use linear regression for?

"linear in the parameters"

$$h(x_i) = \underbrace{w_0 + w_1 \sin(x_i)}_{\text{good:}} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ \sin(x_i) \end{bmatrix}$$

$$h(x_i) = \underbrace{w_0 + w_1 \sin(w_2 x_i + w_3)}_{\text{bad:}}$$

bad: has w 's inside of \sin !

need this model format:

$$\sum w_j \cdot \boxed{\text{no } w\text{'s}}^{(j)}$$

$$h(\vec{x}_i) = w_0 + w_1 x_i^{(1)} + w_2 \frac{x_i^{(1)}}{\log x_i^{(2)}} + w_3 \frac{x_i^{(12)} x_i^{(1)} x_i^{(19)}}{72}$$

this is legal!

$$= \vec{w} \cdot \text{Aug}(\vec{x}_i)$$

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

↳ could map Mon → 1, Tue → 2, ...
but this emphasizes later days
more than earlier days

Solution: one hot encoding

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

5 unique values

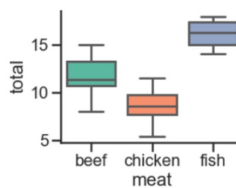
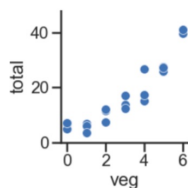
Issue: $Mon + Tue + \dots + Fri =$
Solution: drop one column



Activity 2: Chicken, Beef, or Fish?

Every week, Lauren goes to her local grocery store and buys exactly one pound of meat (either beef, fish, or chicken) but varying amounts of vegetables. We've collected a dataset containing the pounds of vegetables bought, the type of meat bought, and the total bill. Below we display the first few rows of the dataset and two plots generated using the entire (training) dataset.

veg	meat	total
1	beef	13
3	fish	19
2	beef	16
0	chicken	9



Lab 8

In each part below, we provide you with a model that predicts `total` (her total grocery bill), fit to the dataset by minimizing mean squared error. For each model, determine whether **each optimal parameter w^* is positive, negative or exactly 0**. For example, in part (iv), you'll need to provide 3 answers: one for w_0^* , one for w_1^* , and one for w_2^* .

1 $h(\vec{x}_i) = w_0$

2 $h(\vec{x}_i) = w_0 + w_1 \cdot \text{veg}_i$

3 $h(\vec{x}_i) = w_0 + w_1 \cdot \text{meat}=\text{chicken}_i$ (one hot encoded feature for chicken)

4 $h(\vec{x}_i) = w_0 + w_1 \cdot \text{meat}=\text{beef}_i + w_2 \cdot \text{meat}=\text{chicken}_i$

5 $h(\vec{x}_i) = w_0 + w_1 \cdot \text{meat}=\text{beef}_i + w_2 \cdot \text{meat}=\text{chicken}_i + w_3 \cdot \text{meat}=\text{fish}_i$

← Ch. 7

Ch. 8 →

So far, we've minimized

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

using linear algebra and orthogonality.

→ Could we use calculus?

$$R_{sq} : \underbrace{\mathbb{R}^{d+1}}_{\text{domain}} \rightarrow \underbrace{\mathbb{R}}_{\text{codomain}}$$

"vector-to-scalar"
function

Vector-to-scalar function example

$$f(\vec{x}) = x_1^2 + x_2^2 - 3x_1x_2$$

$f: \mathbb{R}^2 \rightarrow \mathbb{R}$ (vector-to-scalar)

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

$$\frac{\partial f}{\partial x_1} = 2x_1 - 3x_2$$

$$\frac{\partial f}{\partial x_2} = 2x_2 - 3x_1$$

New idea: package all partial derivatives of $f(\vec{x})$ into a vector

$$\underbrace{\nabla f(\vec{x})}_{\text{gradient of } f(\vec{x})} = \begin{bmatrix} 2x_1 - 3x_2 \\ 2x_2 - 3x_1 \end{bmatrix}$$

$$\nabla f(\vec{x}) = \begin{bmatrix} 2x_1 - 3x_2 \\ 2x_2 - 3x_1 \end{bmatrix}$$

e.g. if $\vec{x} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$

$$\nabla f\left(\begin{bmatrix} 2 \\ 1 \end{bmatrix}\right) = \begin{bmatrix} 1 \\ -4 \end{bmatrix}$$

The gradient, $\nabla f(\vec{x})$, is a
vector - to - vector function!

Gradient vector

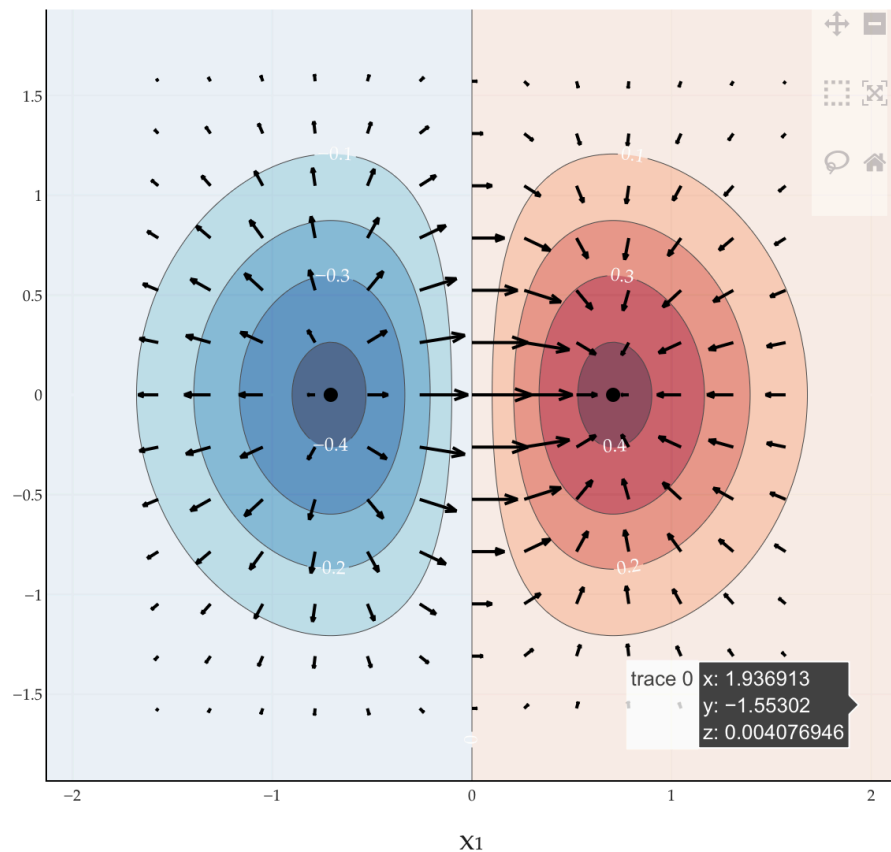
Suppose $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a vector-to-scalar function.

The gradient of f at input \vec{x} is

$$\nabla f(\vec{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

~ describes the direction of steepest ascent!

Gradient field of $f(\vec{x}) = x_1 e^{-x_1^2 - x_2^2}$



Ideally,
to minimize
 $f(\vec{x})$, set
 $\nabla f(\vec{x}) = \vec{0}$

"Big 3 Rules" of Chapter 8.2

(1)

$$f(\vec{x}) = \vec{a} \cdot \vec{x}$$

$$= a_1 x_1 + a_2 x_2 + \dots + a_n x_n$$

where

$$\vec{a}, \vec{x} \in \mathbb{R}^n$$
$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}, \vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\nabla f(\vec{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \vec{a}$$

$$\vec{a} \cdot \vec{x} = \vec{x}^T \vec{a}$$

$$\textcircled{2} \quad \vec{x} \in \mathbb{R}^n$$

$$f(\vec{x}) = \|\vec{x}\|^2 = x_1^2 + x_2^2 + \dots + x_n^2$$

$$\frac{\partial f}{\partial x_i} = 2x_i$$

$$\nabla f(\vec{x}) = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_n \end{bmatrix} = 2\vec{x}$$

Rule ③

"quadratic form"

$$f(\vec{x}) = \vec{x}^T A \vec{x}$$

$$\vec{x} \in \mathbb{R}^n$$

A : $n \times n$ matrix

$$\nabla f(\vec{x}) = (A + A^T) \vec{x}$$

\Rightarrow special case: what if A is symmetric?
 $A = A^T$, so then

$$f(\vec{x}) = \vec{x}^T A \vec{x} \Rightarrow \nabla f(\vec{x}) = 2A \vec{x}$$

\Rightarrow see Lab 8