Lecture 26, 27, 28 - Review, Conclusion



DSC 40A, Fall 2022 @ UC San Diego Mahdi Soleymani, with help from many others

Announcements

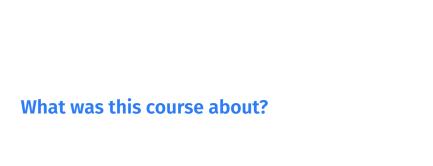
- Homework 8 is due Tuesday Dec. 6. (optional)
- A recording of Discussion 8 (probability review) is posted on the course website and on Campuswire.
- Fill out CAPEs survey.
 - Deadline: Saturday at 8am.
- ► The Final Exam is on **Saturday 12/4 from 7:00PM-10:00PM**.
 - Bring a cheat sheet.
 - Bring a calculator. No other electronic devices are allowed.
 - UCSD ID is required!

Final preparation

- Review the solutions to previous homeworks and groupworks.
 - All except Homework 8 are up.
- ▶ Identify which concepts are still iffy. Re-watch lecture, post on Campuswire, come to office hours.
 - We have many office hours between now and the exam.
- Look at the past exams at https://dsc4oa.com/resources.
 - Watch the probability review discussion.
- Study in groups.
- Make a "cheat sheet".

Agenda

- ► High-level summary of the course.
- ► Review problems.
- ► Conclusion.



Part 1: Supervised learning

The "learning from data" recipe to make predictions:

- 1. Choose a prediction rule. We've seen a few:
 - ightharpoonup Constant: H(x) = h.
 - Simple linear: $H(x) = w_0 + w_1 x$.
 - Multiple linear: $H(x) = w_0 + w_1 x^{(1)} + w_2 x^{(2)} + ... + w_d x^{(d)}$.
- 2. Choose a loss function.
 - Absolute loss: L(h, y) = |y h|.
 - ► Squared loss: $L(h, y) = (y h)^2$.
 - 0-1 loss, UCSD loss, etc.
- 3. Minimize empirical risk to find optimal parameters.
 - Algebraic arguments.
 - Calculus (including vector calculus).
 - Gradient descent.

Part 1: Unsupervised learning

- When learning how to fit prediction rules, we were performing supervised machine learning.
- ▶ We discussed k-Means Clustering, an unsupervised machine learning method.
 - Supervised learning: there is a "right answer" that we are trying to predict.
 - Unsupervised learning: there is no right answer, instead we're trying to find patterns in the structure of the data.

Part 2: Probability fundamentals

- If all outcomes in the sample space S are equally likely, then $P(A) = \frac{|A|}{|S|}$.
- $ightharpoonup \bar{A}$ is the **complement** of event A. $P(\bar{A}) = 1 P(A)$.
- Two events A, B are mutually exclusive if they share no outcomes, i.e. they don't overlap. In this case, the probability that A happens or B happens is $P(A \cup B) = P(A) + P(B)$.
- More generally, for any two events, $P(A \cup B) = P(A) + P(B) P(A \cap B)$.
- The probability that events A and B both happen is $P(A \cap B) = P(A)P(B|A)$.
 - P(B|A) is the probability that B happens given that you know A happened.
 - Through re-arranging, we see that $P(B|A) = \frac{P(A \cap B)}{P(A)}$.

Part 2: Combinatorics

- A sequence is obtained by selecting *k* elements from a group of *n* possible elements with replacement, such that order matters.
 - Number of sequences: n^k .
- A permutation is obtained by selecting *k* elements from a group of *n* possible elements without replacement, such that order matters.
 - Number of permutations: $P(n, k) = \frac{n!}{(n-k)!}$.
- A **combination** is obtained by selecting *k* elements from a group of *n* possible elements without replacement, such that order does not matter.
 - Number of combinations: $\binom{n}{k} = \frac{n!}{(n-k)!k!}$.

Part 2: The law of total probability and Bayes' theorem

- A set of events $E_1, E_2, ..., E_k$ is a partition of S if each outcome in S is in exactly one E_i .
- The law of total probability states that if A is an event and $E_1, E_2, ..., E_k$ is a partition of S, then

$$P(A) = P(E_1) \cdot P(A|E_1) + P(E_2) \cdot P(A|E_2) + \dots + P(E_k) \cdot P(A|E_k)$$

$$= \sum_{i=1}^{k} P(E_i) \cdot P(A|E_i)$$

Bayes' theorem states that

$$P(B|A) = \frac{P(B) \cdot P(A|B)}{P(A)}$$

We often re-write the denominator P(A) in Bayes' theorem using the law of total probability.

Part 2: Independence and conditional independence

- Two events A and B are independent when knowledge of one event does not change the probability of the other event.
 - Equivalent conditions: P(B|A) = P(B), P(A|B) = P(A), $P(A \cap B) = P(A) \cdot P(B)$.
- Two events A and B are conditionally independent if they are independent given knowledge of a third event, C.
 - ► Condition: $P((A \cap B)|C) = P(A|C) \cdot P(B|C)$.
- In general, there is no relationship between independence and conditional independence.
- See pinned post on Campuswire for clarification.

Part 2: Naive Bayes

- In classification, our goal is to predict a discrete category, called a **class**, given some features.
- ► The Naive Bayes classifier works by estimating the numerator of *P*(class|features) for all possible classes.
- It uses Bayes' theorem:

$$P(\text{class}|\text{features}) = \frac{P(\text{class}) \cdot P(\text{features}|\text{class})}{P(\text{features})}$$

► It also uses a "naive" simplifying assumption, that features are conditionally independent given a class:

$$P(\text{feature}_1|\text{class}) \cdot P(\text{feature}_2|\text{class}) \cdot \dots$$

Skipped problems

Review problems

Example: Clustering and combinatorics

- Suppose we have a dataset of 15 points, each with two features (x_1, x_2) . In the dataset, there exist 3 "natural" clusters, each of which contain 5 data points.
- Recall that in the k-Means Clustering algorithm, we initialize k centroids by choosing k points at random from our dataset. Suppose k = 3.

1.	What's the probability that all three initial centroids are initialized in the same natural cluster?
2.	What's the probability that all three initial centroids are initialized in different natural clusters?

Example: basketball

Suppose we have 6 basketball players who want to organize themselves into 3 basketball teams of 2 players each. Suppose

we have three teams, "Team USA", "Team China", and "Team Lithuania". How many ways can these teams be formed?

Example: basketball, again

Suppose we have 6 basketball players who want to organize themselves into 3 basketball teams of 2 players each. Now, suppose the teams are irrelevant, and all we care about is the unique pairings themselves. How many ways can these 6 players be split into 3 teams?

Example: high school

A certain high school has 80 students: 20 freshmen, 20 sophomores, 20 juniors, and 20 seniors. If a random sample of 20 students is drawn without replacement, what is the probability that the sample contains 5 students in each grade level?

Example: high school, again

A certain high school has 80 students: 20 freshmen, 20 sophomores, 20 juniors, and 20 seniors. If a random sample of 20 students is drawn with replacement, what is the probability that all students in the sample are from the same grade level?

Example: bitstrings

What is the probability of a randomly generated bitstring of length 5 having the same first two bits? Assume that each bit is equally likely to be a 0 or a 1.

Example: bitstrings, again

What is the probability of a randomly generated bitstring of length 5 having the same first two bits, if we know that the bitstring has exactly four 0s? Assume that each bit is equally likely to be a 0 or a 1.

Example: Two-sided cards

Source

Suppose we have 3 cards identical in form except that both sides of the first card are colored red, both sides of the second card are colored black, and one side of the third card is colored red and the other side is colored black. The 3 cards are mixed up in a hat, and 1 card is randomly selected and put down on the ground. If the upper side of the chosen card is colored red, what is the probability that the other side is colored black?

Marbles

Suppose you are given two jars. Jar I contains one black and 4 white marbles, and Jar II contains 4 black and 6 white marbles. If a jar is selected at random and a marble is chosen,

- What is the probability that the marble chosen is a black marble?
- ► If the chosen marble is black, what is the probability that it came from Jar I?
- If the chosen marble is black, what is the probability that it came from Jar II?
 Source

Given a 2 x 2 matrix:

$$X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix},$$

its inverse (if exists) is given by:

$$X^{-1} = \frac{1}{\det(X)} \cdot \begin{pmatrix} x_{22} & -x_{12} \\ -x_{21} & x_{11} \end{pmatrix},$$

where det(X) is the determinant of X:

$$\det(X) = x_{11} x_{22} - x_{21} x_{12}.$$

Example

$$det(X) = (2x+1) - 3x1 = -1$$

$$X = \begin{pmatrix} 2 & 3 \\ 1 & +1 \end{pmatrix}$$

what is X^{-1} ?

what is
$$X^{-1}$$
?
$$\frac{1}{-1} \begin{pmatrix} 1 & -3 \\ -1 & 2 \end{pmatrix} = \begin{pmatrix} -1 & 3 \\ 1 & -2 \end{pmatrix} = \begin{pmatrix} -1 & 3 \\ 1 & -2 \end{pmatrix}$$

$$X = \begin{pmatrix} 2 & 3 \\ 1 & +1 \end{pmatrix},$$

Example

Example
$$de+(X) = 6x+1 - 3x2 = 6-6=0$$

=> X has no inverse.

Given: $X = \begin{pmatrix} 6 & 3 \\ 2 & +1 \end{pmatrix}$

what is X^{-1} ?

Linear regression

Data: (3,2),(5,1), and (4,3).

$$X = \begin{pmatrix} 1 & 3 \\ 1 & 5 \\ 1 & 4 \end{pmatrix}$$
The design matrix

$$\overrightarrow{y} = \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the sign matrix

$$\overrightarrow{y} = \begin{pmatrix} 3 \\ 1 \\ 3 \end{pmatrix}$$
The property of the si

 $\chi^{T}\chi = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & 5/124 & 50/ & 4 \end{pmatrix} = \begin{pmatrix} 1 & \sqrt{31} & 12/\\ 3 & \sqrt{31} & 12$









$$det(x^{T}x) = 3x60 - 12x12 = 6$$

$$(x^{T}x)^{-1} = \frac{1}{6} \begin{pmatrix} 50 & -12 \\ -12 & 3 \end{pmatrix}$$

$$(x^{T}x)^{-1}x^{T} = \frac{1}{6} \begin{pmatrix} 50 & -12 \\ -12 & 3 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 3 & 5 & 4 \end{pmatrix}$$

$$= \frac{1}{6} \begin{pmatrix} 14 & 38 & 2 \\ -3 & 3 & 0 \end{pmatrix}$$

$$= (x^{T}x)^{-1}x^{T} \vec{x} = \frac{1}{6} \begin{pmatrix} -3 & 3 & 0 \\ -1/2 & 4 & 4 \end{pmatrix}$$

$$= \frac{1}{6} \begin{pmatrix} 14 & 38 & 2 \\ -1/2 & 4 & 4 \end{pmatrix}$$

$$= \frac{1}{6} \begin{pmatrix} 14 & 38 & 2 \\ -1/2 & 4 & 4 \end{pmatrix}$$

oid
$$(e^{x})'=e^{x}$$

$$(e^{x})'=e^{x}$$

$$(e^{x})'=e^{x}$$

$$(e^{x})'=e^{x}$$

We have the Sigmoid function: $\sigma(x) = \frac{1}{1+e^{-x}}$. What is its first

derivative?
$$\sigma'(x) = \frac{1}{2} \sigma(x) = \frac{1}{2} (1 + e^{-x})^{-1} = \frac{1}{2} (1 + e^{-x})^$$

$$\sigma'(x) = \frac{1}{3x} \sigma(x) = \frac{1}{3x} (1 + e^{-x})^{-1} = \frac{1}{3x} (1 + e^{-x})^{-1}$$

$$(-1) \left(1 + e^{-x}\right)^{-2} \frac{d}{dx} \left(1 + e^{-x}\right) = \frac{-x}{(1 + e^{-x})^{-2}} \left(-1\right) e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^{-2}} = 0 \text{ (1)} \left(1 - 0 \text{ (1)}\right)$$

$$(-1) \left(1 + e^{-x}\right)^{-2} \frac{d}{dx} \left(1 + e^{-x}\right) = \frac{e^{-x}}{(1 + e^{-x})^{2}}$$

$$(-1) \left(1 + e^{-x}\right)^{-2} \left(-1\right) e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^{2}}$$

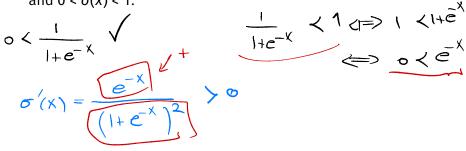
$$(-1) \left(1 + e^{-x}\right)^{-2} \left(-1\right) e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^{2}}$$

$$(-1) \left(1 + e^{-x}\right)^{-2} \left(-1\right) e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^{2}}$$

$$(-1) \left(1 + e^{-x}\right)^{-2} \left(-1\right) e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^{2}}$$

Sigmoid

Show that the Sigmoid function is monotonically increasing and $0 < \sigma(x) < 1$.



Logistic Regression

Given a hypothesis

$$h(\vec{x}; \vec{w}, w_0) = \sigma(\vec{x} \cdot \vec{w} + w_0).$$

In logistic regression, we have:

$$p(y = 1 | \vec{x}; \vec{w}, w_0) = \sigma(\vec{x} \cdot \vec{w} + w_0).$$

Let's derive the partial derivative of h with respect to \vec{w} and $w_0!$

$$\frac{Jh}{\partial W_0} = \sigma'(z) \cdot \frac{Jz}{JW_0}$$

$$\frac{\partial h}{\partial \overline{W}} = \sigma'(z) \cdot \frac{\partial z}{\partial \overline{W}}$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = 1$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{W} + W_0 \end{array} \right) = X$$

$$\frac{d}{dW_0} \left(\begin{array}{c} \overrightarrow{X} \cdot \overrightarrow{$$

 $\frac{\partial h}{\partial \vec{w}} = \sigma(\vec{z}) \cdot \frac{\partial \vec{z}}{\partial \vec{w}} = \sigma(\vec{z}) \left(1 - \sigma(\vec{z}) \cdot \vec{x}\right)$ $\text{fint} = \sigma(\vec{x} \cdot \vec{w} + w_0) \left(1 - \sigma(\vec{x} \cdot \vec{w} + w_0)\right) 1$ $\text{Second} \qquad \sigma(\vec{x} \cdot \vec{w} + w_0) \left(1 - \sigma(\vec{x} \cdot \vec{w} + w_0)\right) \vec{x}$

Gradient

What is the gradient of
$$f(x,y) = x^2 + \sin(y+x) + xe^{-y} + y$$

what is the gradient of
$$f(x,y) = x^2 + \sin(y + x) + xe^{-y} + y$$

$$\frac{\partial f}{\partial x} = 2 \times + (1) G_S(y + x) + e^{-y} + 0$$

$$\frac{\partial f}{\partial x} = 0 + (1) G_S(y + x) + \chi(-e^{-y}) + 1$$

$$\nabla f(0,0) = \begin{pmatrix} 6 + 1 + 1 + 0 \\ 0 + 1 + 0 + 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

Conclusion

Learning objectives

At the start of the quarter, we told you that by the end of DSC 40A, you'll...

- understand the basic principles underlying almost every machine learning and data science method.
- be better prepared for the math in upper division: vector calculus, linear algebra, and probability.
- be able to tackle problems such as:
 - How do we know if an avocado is going to be ripe before we eat it?
 - How do we teach a computer to read handwritten text?
 - How do we predict a future data scientist's salary?

What's next?

In DSC 40A, we just scratched the surface of the theory behind data science. In future courses, you'll build upon your knowledge from DSC 40A, and will learn:

- More supervised learning.
 - Logistic regression, decision trees, neural networks, etc.
- More unsupervised learning.
 - Other clustering techniques, PCA, etc.
- More probability.
 - Random variables, distributions, etc.
- More connections between all of these areas.
 - For instance, you'll learn how probability is related to linear regression.
- More practical tools.

Thank you!

- The other instructor, Dr. Truong Son Hy.
- This course would not have been possible without our TA: Pushkar Bhuse.
- It also would not have been possible without our 8 tutors: Yuxin Guo, Weiyue(Larry) Li, Vivian Lin, Karthikeya Manchala, Shiv Sakthivel, Aryaman Sinha, Jessica Song and Yujia(Joy) Wang.
- You can contact them with any questions at dsc40a.com/staff.

