# CS 260: Foundations of Data Science

Prof. Thao Nguyen

Fall 2024

# Admin

- **Lab 2** grades & feedback will be posted on Wednesday

- **Lab 3** due tonight

- **Lab 4** posted, due next Monday at midnight

- **Lecture Notes**

# Peer Tutoring

- **Student tutors** (Fejiro Anigbro, Darshan Mehta)

- **Flexible hours**

- **Free!**

# TECH TALKS 2024

## OCTOBER 7, 8 & 9TH | 6-8PM EST

*Sign up for a 30 minute virtual informational interview with a Tri-Co alum to gain tech career insights!*

*Alumni will represent various tech roles including software engineering and development, data science, tech consulting, product management and biotech.*

| OCT 7 | OCT 8 | OCT 9 |
|---|---|---|
| Accenture | Bristol Myers Squibb | The Walt Disney Company |
| FERMAT Commerce | Community.com | Fresh Tracks Insights |
| Grubhub | C3 Presents (Live Nation) | Meta |
|  | Opower (Oracle) | Grubhub |

TRI-COLLEGE RECRUITING CONSORTIUM

HAVERFORD · BRYN MAWR · SWARTHMORE

# Outline for today

- Recap SGD (stochastic gradient descent)

- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation

- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Outline for today

- Recap SGD (stochastic gradient descent)

- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation

- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Stochastic Gradient Descent for Linear Regression

Key Idea: take the derivative of **one datapoint** at a time and use that to update w

# Stochastic Gradient Descent for Linear Regression

SGD for Linear Regression

Start with $\vec{w} = \vec{0}$ (vector of zeros)

~~for~~ while (epoch) iteration $t$:                    not for Lab 3

for $i = 1, 2 \cdots n$: (shuffle)

$$\vec{w} \leftarrow \vec{w} - \alpha \left( \vec{w} \cdot \vec{x}_i - y_i \right) \vec{x}_i$$

Step size        derivative

if check convergence

$$\left| J(\vec{w}^t) - J(\vec{w}^{t+1}) \right| < \varepsilon \qquad \varepsilon \text{ is very small}$$

$\Rightarrow$ Stop!

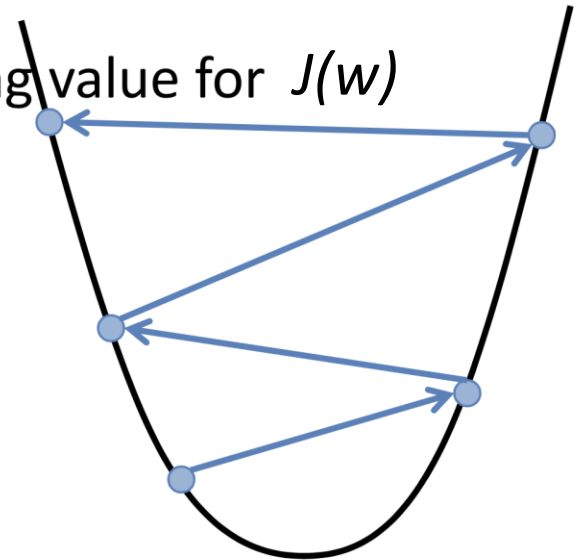# Choosing the step size alpha

$\alpha$ **too small**

slow convergence

$\alpha$ **too large**

increasing value for $J(w)$

- may overshoot minimum
- may fail to converge (may even diverge)

# Pros and Cons

(Analytic Solution)

## Gradient Descent

- requires multiple iterations
- need to choose $\alpha$
- works well when $p$ is large
- can support online learning

## Normal Equations

- non-iterative
- no need for $\alpha$
- slow if $p$ is large
  - matrix inversion is $O(p^3)$

# Outline for today

- Recap SGD (stochastic gradient descent)


- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation


- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Binary classification examples

- Transactions that indicate credit card fraud

- Accounts that are bots

- Detecting which scans show tumors

- Prenatal test for Down's Syndrome

- Finding genes under natural selection

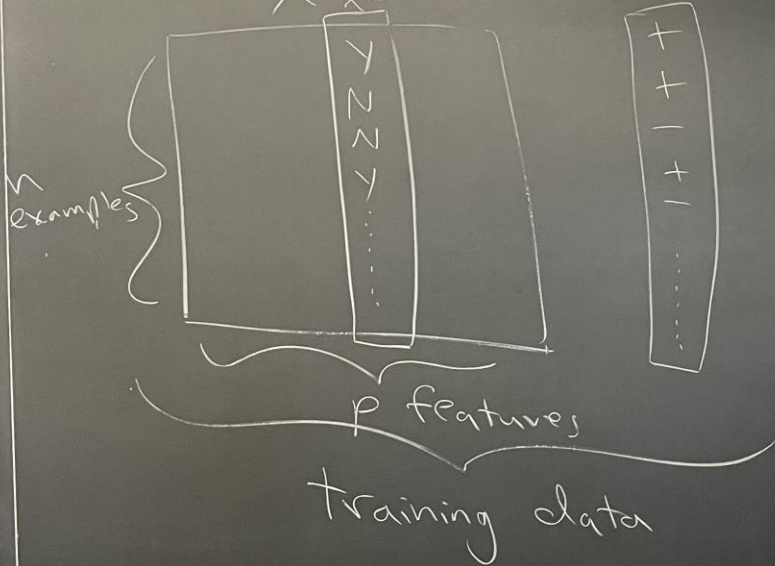- Regions of the environment that contains the object the robot is searching for

In all these examples, we are trying to find unusual items ("needle in a haystack") -- we call these *positives*

<u>Classification</u>

X fever     y (disease)

$\overbrace{\phantom{xxxx}}^{n\ examples}$   p features

training data

$+ \Rightarrow$ disease
$- \Rightarrow$ no disease

<u>Model</u>: decision tree with a single feature ("stump")

fever

Y      N

$P_{pos} = \dfrac{6}{8}$     n=15     $P_{pos} = \dfrac{3}{7}$

$P_{pos} =$ prob of positive
Y

new idea : use probabilities
to classify      test examples

$$\vec{x}_{test} = [ \ \cdots \ \overset{fever}{N} \ \cdots \ ]^T$$

threshold 0.5 $\Rightarrow$ $\boxed{\hat{y}_{test} = \ominus}$ no disease

threshold 0.25 $\Rightarrow$ $\boxed{\hat{y}_{test} = \oplus}$ disease

$\boxed{P_{pos} \geqslant threshold \Rightarrow classify \ \oplus}$

# Outline for today

- Recap SGD (stochastic gradient descent)

- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation

- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Goals of Evaluation

- Think about what metrics are important for the problem at hand

- Compare different methods or models on the same problem

- Common set of tools that other researchers/users can understand

# Training and Testing
## (high-level idea)

- **Separate** data into "**train**" and "**test**"
  - *n* = num training examples
  - *m* = num testing examples

- **Fit** (create) the model using **training data**
  - e.g. sea_ice_1979-2012.csv

- **Evaluate** the model using **testing data**
  - e.g. sea_ice_2013-2020.csv

$$\frac{65+13}{100} = 78\%$$

pred

|  | − | + |
|---|---|---|
| truth − | 65 | 15 |
| + | 7 | 13 |

$$accuracy =$$

## test data

$$m = 100$$

thresh = 0.5

| 50 | 30 |
|---|---|
| 1 | 19 |

thresh 0.25

acc = 69%

$$= \frac{\# \, correct}{m}$$

$$= \boxed{\frac{1}{m} \sum_{i=1}^{m} \mathbb{1}(\hat{y}_i == y_i)}$$

80 negatives
20 positives

| 76 | 4 |
|---|---|
| 11 | 9 |

thresh 0.75

Note: all the same model, different thresholds!

# Confusion Matrices

Predicted class

|  | Negative | Positive |
|---|---|---|
| **Negative** | True negative (TN) | False positive (FP) |
| **Positive** | False negative (FN) | True positive (TP) |

True class

# Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N (total number of true negatives) |
| Positive | False negative (FN) "miss" | True positive (TP) | P (total number of true positives) |
|  | N* (what we said was negative) | P* (what we said was positive "flagged") |  |

True class

# Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| **Negative** | True negative (TN) ✔ | False positive (FP) "false alarm" ✖ | N |
| **Positive** | False negative (FN) "miss" ✖ | True positive (TP) ✔ | P |
|  | N* | P* |  |

True class

# Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Error:

(FN+FP)/(TN+FP+FN+TP)

= (FN+FP)/(N+P)

# Confusion Matrices

Predicted class

|  | Negative | Positive | |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
| | N* | P* | |

True class

Accuracy = 1-Error:

(TN+TP)/(TN+FP+FN+TP)

= (TN+TP)/(N+P)

# Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Precision:

$$TP/(FP+TP) = TP/P*$$

# Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Recall
(True Positive Rate):

$TP/(FN+TP) = TP/P$

# Confusion Matrices



False Positive Rate:

FP/(TN+FP) = FP/N

# Precision and Recall

- <u>Precision</u>: of all the "flagged" examples, which ones are actually relevant (i.e. positive)?

  (Purity)

- <u>Recall</u>: of all the relevant results, which ones did I actually return?

  (Completeness)