# CS 260: Foundations of Data Science

Prof. Thao Nguyen

Fall 2024

HAVERFORD
COLLEGE

# 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 |

# 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



$$\nabla J = \begin{bmatrix} \dfrac{\partial J}{\partial w_0} \\ \vdots \\ \dfrac{\partial J}{\partial w_p} \end{bmatrix}$$

Handout 6

1 + 2

derivative w.r.t. one example

$$J(\vec{w}) = \frac{1}{2} \sum_{i=1}^{n} \left( \underbrace{\vec{w} \cdot \vec{x}_i}_{pred} - \underbrace{y_i}_{truth} \right)^2$$

derivative is very large & unstable

$$\nabla J(\vec{w})_{\vec{x}_i} = \left( \vec{w} \cdot \vec{x}_i - y_i \right) \vec{x}_i$$

datapoint     scalar     vector

$$\vec{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}$$

# Stochastic Gradient Descent for Linear Regression



SGD) for Linear

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)

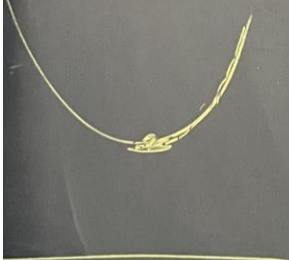$$\vec{w} \leftarrow \vec{w} - \alpha (\vec{w} \cdot \vec{x}_i - y_i) \vec{x}_i$$

↑ step size    derivative

if check convergence

$$\left| J(\vec{w}^t) - J(\vec{w}^{t+1}) \right| < \xi \quad \leftarrow \quad \xi \text{ is very small}$$
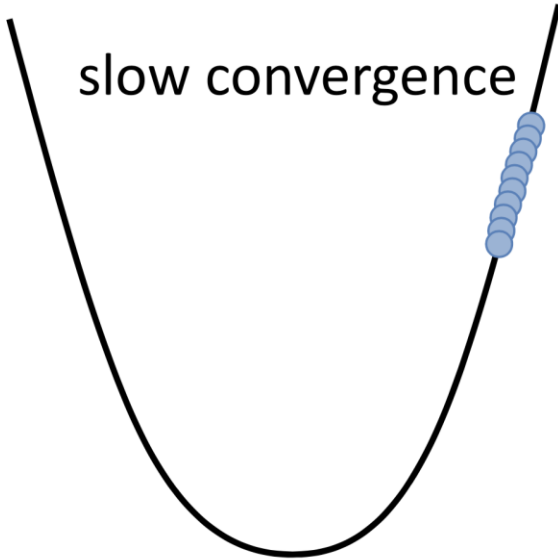
⟹ Stop!

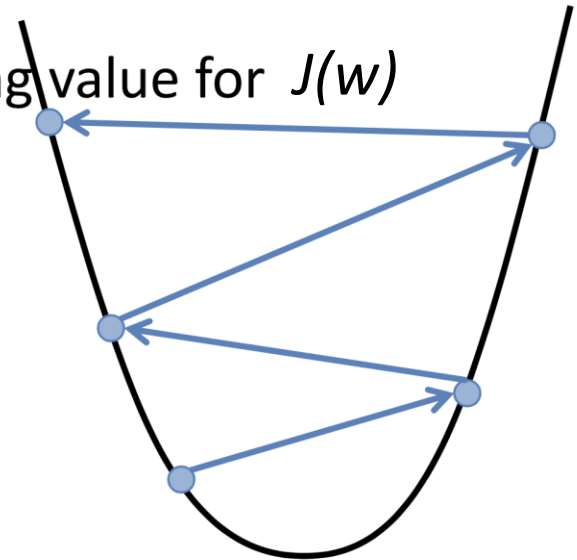# Choosing the step size alpha



$\alpha$ 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

- Finding regions of the genome with high recombination rate ("hotspots")

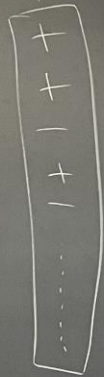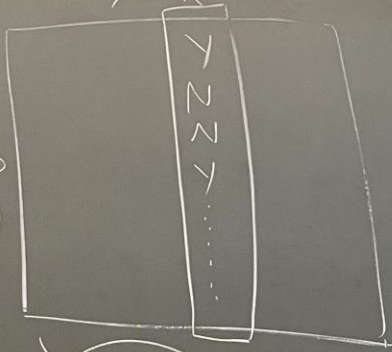In all these examples, we are trying to find unusual items ("needle in a haystack") -- we call these *positives*

Classification

$X$ fever     $y$ (disease)

n examples

p features

training data

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

Model: decision tree with a single feature ("stump")

fever

Y    N

$P_{pos} = \frac{6}{8}$

$n = 15$

$P_{pos} = \text{prob of positive} = Y$

$P_{pos} = \frac{3}{7}$

new idea : use probabilities

to classify    test examples

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

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

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

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

# Handout 7

# 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 |

$$\text{accuracy} = \frac{\text{\# correct}}{m}$$

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

Note: all the same model, different thresholds!

test data

$m = 100$

thresh = 0.5

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

thresh 0.25

$$\text{acc} = 69\%$$

80 negatives
20 positives

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

thresh 0.75

# Confusion Matrices

Predicted class

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

True class