

10-607 Computational Foundations for Machine Learning

Perceptron Mistake Bound & Computational Complexity

Instructor: Pat Virtue

Plan

Perceptron Algorithm

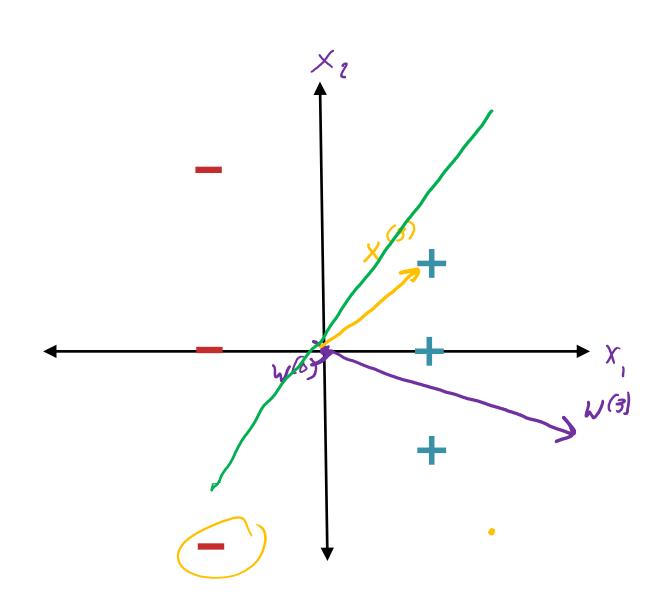
Perceptron Mistake Bound Theory

- Background: projections, distances, and margin
- Proof of mistake bound as an example application

Computational Complexity

- How fast is your code/algorithm?
- Counting operations
- Big-O
- Complexity classes

Sketch of algorithm Initialize w = 0Loop Given a point \mathbf{x} predict $\hat{y} = sign(\mathbf{w}^T \mathbf{x})$ Given actual label *y*: If $y \neq \hat{y}$ If y = -1, $w \leftarrow w - x$



Perceptron Algorithm: Example

Example:
$$(-1,2) - \times$$

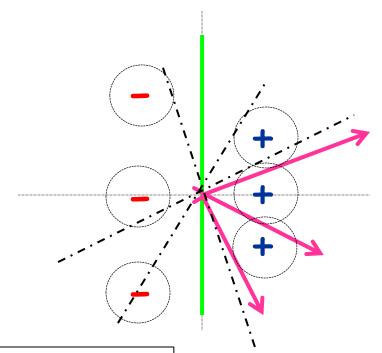
$$(1,0) + \checkmark$$

$$(1,1) + \times$$

$$(-1,0) - \checkmark$$

$$(-1,-2) - \times$$

$$(1,-1) + \checkmark$$



Perceptron Algorithm: (without the bias term)

- Set t=1, start with all-zeroes weight vector w_1 .
- Given example x, predict positive iff $w_t \cdot x \ge 0$.
- On a mistake, update as follows:
 - Mistake on positive, update $w_{t+1} \leftarrow w_t + x$
 - Mistake on negative, update $w_{t+1} \leftarrow w_t x$

$$w_{1} = (0,0)$$

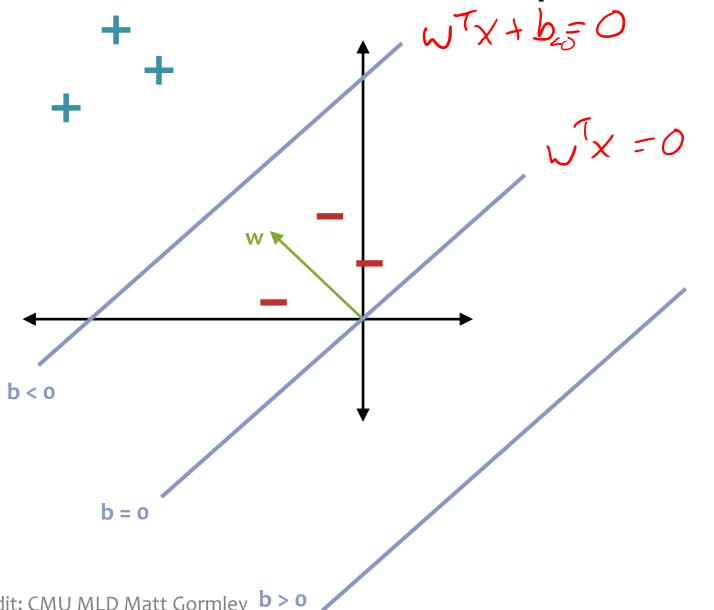
$$w_{2} = w_{1} - (-1,2) = (1,-2)$$

$$w_{3} = w_{2} + (1,1) = (2,-1)$$

$$w_{4} = w_{3} - (-1,-2) = (3,1)$$

$$w_{4} = -(-1,2) + (1,1) - (-1,-2)$$

Intercept Term



Q: Why do we need an intercept term?

A: It shifts the decision boundary off the origin

Q: What should happen to b during the perceptron algorithm

A: Two cases

- 1. Increasing b shifts the decision boundary towards the negative side
- 2. Decreasing b shifts the decision boundary towards the positive side

Sketch of algorithm

```
Initialize w = 0, b = 0

Loop

Given a point x predict \hat{y} = sign(w^Tx)

Given actual label y:

If true y \neq \hat{y}

If true y = +1, w \leftarrow w + x, b = b + 1

If true y = -1, w \leftarrow w - x, b = b - 1
```

$x_{\text{orig}} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \longrightarrow \bar{X} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$

Sketch of algorithm

Initialize
$$\theta = 0$$

Loop

Given a point x predict $\hat{y} = sign(\theta^T x)$

Given actual label *y*:

If true
$$y \neq \hat{y}$$

If true $y = +1$, $\theta \leftarrow \theta + x$

If true $y = -1$, $\theta \leftarrow \theta - x$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + yx$$

$$\vec{U} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \qquad \vec{\theta} = \begin{bmatrix} b \\ v_1 \\ v_2 \end{bmatrix}$$

$$\frac{\partial^T x}{\uparrow} = w^T x_{arig} + b$$

Learning for Perceptron if we have a fixed training dataset, D.

Algorithm 1 Perceptron Learning Algorithm

```
1: procedure Perceptron(\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\})
                                                            ⊳ Iı ''' arameters
   	heta \leftarrow 0
2:
3: while not converged do
              for i \in \{1, 2, \dots, N\} do
                                                               ▷ For each example
4:
5:
                   \hat{y} \leftarrow \mathsf{sign}(\boldsymbol{\theta}^T \mathbf{x}^{(i)})
                                                                                ▶ Predict
                  if \hat{y} \neq y^{(i)} then
                                                                            ▶ If mistake
6:
          \theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}

    □ Update parameters

         return \theta
8:
```

Learning for Perceptron if we have a fixed training dataset, D.

Algorithm 1 Perceptron Learning Algorithm

```
1: procedure Perceptron(\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\})
                                                                      ⊳ I · · · · arameters
        	heta \leftarrow 0
2:
    while not converged do
                for i \in \{1, 2, \dots, N\} do
                                                                        ▷ For each example
4:
                      \hat{y} \leftarrow \mathsf{sign}(\boldsymbol{\theta}^T \mathbf{x}^{(i)})
5:
                                                                                              ▶ Predict
                      if \hat{y} \neq y^{(i)} then
                                                                                         ▶ If mistake
6:
                            \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y^{(i)} \mathbf{x}^{(i)}
7:
```

Implementation Trick: same behavior as our "add on positive mistake and subtract on negative mistake" version, because y⁽ⁱ⁾ takes care of the sign

8:

return θ

Plan

Perceptron Algorithm

Perceptron Mistake Bound Theory

- Background: projections, distances, and margin
- Proof of mistake bound as an example application

Computational Complexity

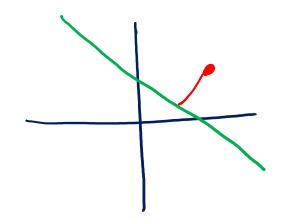
- How fast is your code/algorithm?
- Counting operations
- Big-O
- Complexity classes

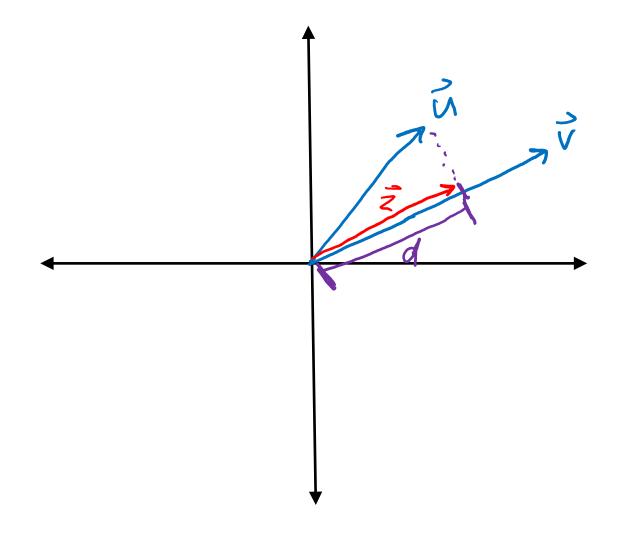
Projection

Projection of ${\boldsymbol u}$ on to ${\boldsymbol v}$

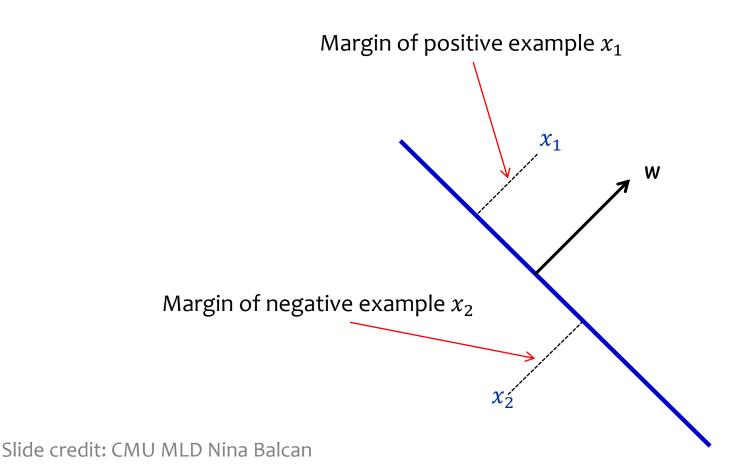
$$= \frac{u^T v}{\|v\|_2}$$

$$\vec{z} = d \frac{\vec{v}}{\|v\|_2}$$



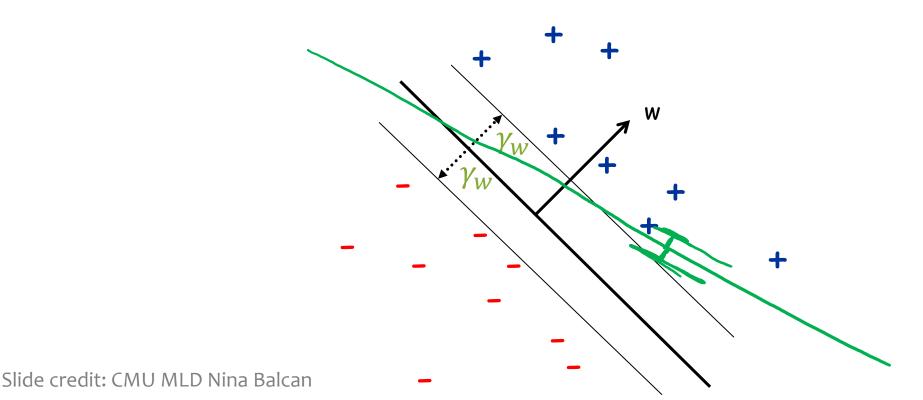


Definition: The margin of example x w.r.t. a linear sep. w is the distance from x to the plane $w \cdot x = 0$ (or the negative if on wrong side)



Definition: The margin of example x w.r.t. a linear sep. w is the distance from x to the plane $w \cdot x = 0$ (or the negative if on wrong side)

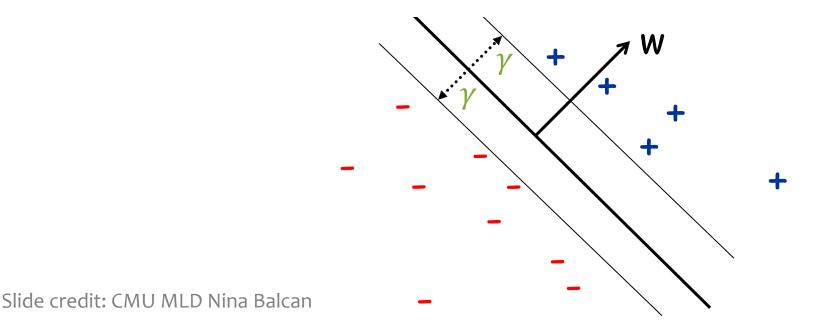
Definition: The margin γ_w of a set of examples S wrt a linear separator w is the smallest margin over points $x \in S$.



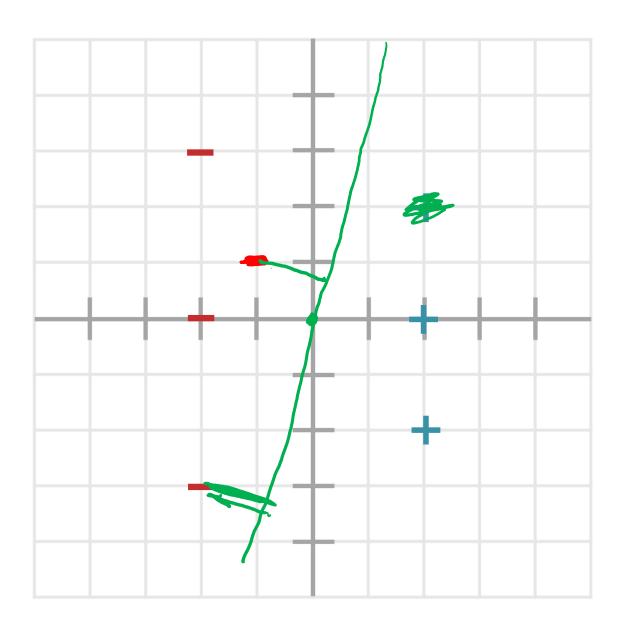
Definition: The margin of example x w.r.t. a linear sep. w is the distance from x to the plane $w \cdot x = 0$ (or the negative if on wrong side)

Definition: The margin γ_w of a set of examples S wrt a linear separator w is the smallest margin over points $x \in S$.

Definition: The margin γ of a set of examples S is the maximum γ_w over all linear separators w.

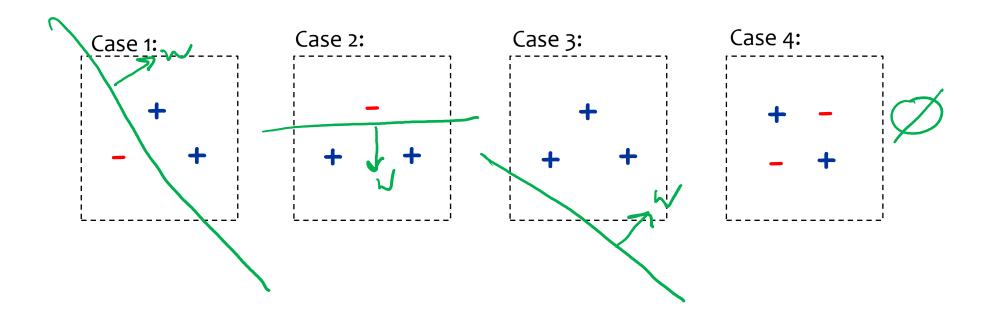


What is the margin for this dataset?



Linear Separability

Def: For a **binary classification** problem, a set of examples *S* is **linearly separable** if there exists a linear decision boundary that can separate the points

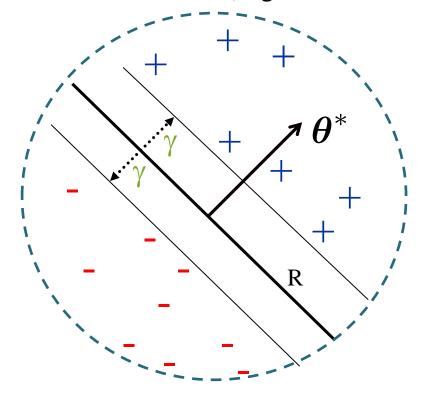


ANALYSIS OF PERCEPTRON

Perceptron Mistake Bound

Guarantee: If data has margin γ and all points inside a ball of radius R, then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes; algo is invariant to scaling.)



Perceptron Mistake Bound

Guarantee: If data has margin γ and all points inside a ball of radius R, then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes; algo is invariant to scaling.)



Def: We say that the perceptron algorithm has **converged** if it stops making mistakes on the training data (perfectly classifies the training data).

Main Takeaway: For linearly separable data, if the perceptron algorithm cycles repeatedly through the data, it will converge in a finite # of steps.

Perceptron Mistake Bound

Theorem 0.1 (Block (1962), Novikoff (1962)).

Given dataset: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$.

Suppose:

1. Finite size inputs: $||x^{(i)}||_{2} \leq R$ 2. Linearly separable data: $\exists \theta^*$ s.t. $||\theta^*||_{2} = 1$ and $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \ge \gamma, \forall i$

Then: The number of mistakes made by the Perceptron

algorithm on this dataset is

$$k \le (R/\gamma)^2$$

Analysis: Percept Misunderstanding:

Perceptron Mistake Boun

Theorem 0.1 (Block (1962), Novikoff (19

Common
Misunderstanding:
The radius is
centered at the
origin, not at the
center of the
points.

- 1. Finite size inputs: $||x^{(i)}|| \leq R$
- 2. Linearly separable data: $\exists \theta^*$ s.t. $||\theta^*|| = 1$ and $y^{(i)}(\theta^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

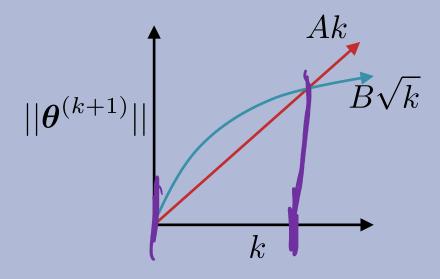
Then: The number of mistakes made by the Perceptron algorithm on this dataset is

$$k \le (R/\gamma)^2$$

Proof of Perceptron Mistake Bound:

We will show that there exist constants A and B s.t.

$$|Ak \le ||\boldsymbol{\theta}^{(k+1)}|| \le B\sqrt{k}$$



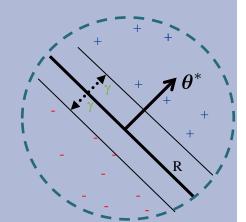
Theorem 0.1 (Block (1962), Novikoff (1962)).

Given dataset: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$.

Suppose:

- 1. Finite size inputs: $||x^{(i)}|| \leq R$
- 2. Linearly separable data: $\exists \boldsymbol{\theta}^*$ s.t. $||\boldsymbol{\theta}^*|| = 1$ and $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \ge \gamma, \forall i$

Then: The number of mistakes made by the Perceptron algorithm on this dataset is



$$k \le (R/\gamma)^2$$

Algorithm 1 Perceptron Learning Algorithm (Online)

```
1: procedure Perceptron(\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \ldots\})
                                                          ▷ Initialize parameters
         \theta \leftarrow \mathbf{0}, k = 1
     for i \in \{1, 2, ...\} do
                                                                          ⊳ For each example
              if y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) < 0 then
                                                                                       ▷ If mistake
                    \boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + y^{(i)} \mathbf{x}^{(i)}
                                                                   k \leftarrow k + 1
         return \theta
7:
```

Slide credit: CMU MLD Matt Gormle

Analysis: Perceptron

Proof of Perceptron Mistake Bound:

Part 1: for some A, $Ak \leq ||\boldsymbol{\theta}^{(k+1)}||$

$$\boldsymbol{\theta}^{(k+1)} \cdot \boldsymbol{\theta}^* = (\boldsymbol{\theta}^{(k)} + y^{(i)} \mathbf{x}^{(i)}) \boldsymbol{\theta}^*$$

by Perceptron algorithm update

$$= \boldsymbol{\theta}^{(k)} \cdot \boldsymbol{\theta}^* + y^{(i)} (\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)})$$

$$\geq \boldsymbol{\theta}^{(k)} \cdot \boldsymbol{\theta}^* + \gamma$$

by assumption

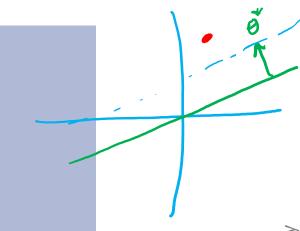
$$\Rightarrow \boldsymbol{\theta}^{(k+1)} \cdot \boldsymbol{\theta}^* \ge k\gamma$$

by induction on k since $\theta^{(1)} = \mathbf{0}$

$$\Rightarrow ||\boldsymbol{\theta}^{(k+1)}|| \geq k\gamma$$

since
$$||\mathbf{w}|| \times ||\mathbf{u}|| \geq \mathbf{w} \cdot \mathbf{u}$$
 and $||\theta^*|| = 1$

Cauchy-Schwartz inequality



Slide credit: CMU MLD Matt Gormley

Analysis: Perceptron

2 (u; + v;)2

Proof of Perceptron Mistake Bound:

Part 2: for some B, $||\boldsymbol{\theta}^{(k+1)}|| \leq B\sqrt{k}$

$$||\boldsymbol{\theta}^{(k+1)}||^2 = ||\boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)}||^2$$

by Perceptron algorithm update

$$= ||\boldsymbol{\theta}^{(k)}||^2 + (y^{(i)})^2||\mathbf{x}^{(i)}||^2 + 2y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)})$$

$$\leq ||\boldsymbol{\theta}^{(k)}||^2 + (y^{(i)})^2 ||\mathbf{x}^{(i)}||^2$$

since kth mistake $\Rightarrow y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$

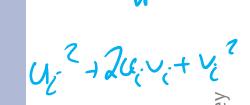
$$= ||\boldsymbol{\theta}^{(k)}||^2 + R^2$$

since $(y^{(i)})^2 ||\mathbf{x}^{(i)}||^2 = ||\mathbf{x}^{(i)}||^2 = R^2$ by assumption and $(y^{(i)})^2 = 1$

$$\Rightarrow ||\boldsymbol{\theta}^{(k+1)}||^2 \le kR^2$$

by induction on k since $(\theta^{(1)})^2 = 0$

$$\Rightarrow ||\boldsymbol{\theta}^{(k+1)}|| \leq \sqrt{k}R$$



Proof of Perceptron Mistake Bound:

Part 3: Combining the bounds finishes the proof.

$$k\gamma \le ||\boldsymbol{\theta}^{(k+1)}|| \le \sqrt{k}R$$
$$\Rightarrow k \le (R/\gamma)^2$$

The total number of mistakes must be less than this



Plan

Perceptron Algorithm

Perceptron Mistake Bound Theory

- Background: projections, distances, and margin
- Proof of mistake bound as an example application

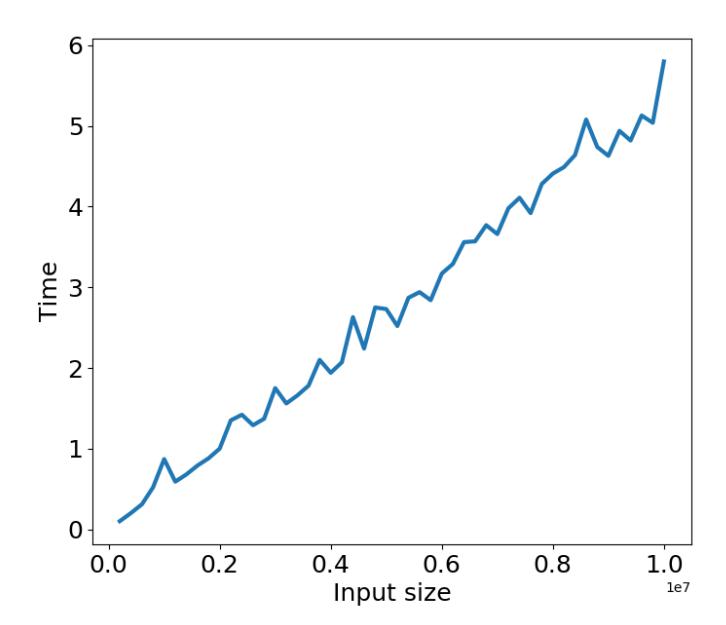
Computational Complexity

- How fast is your code/algorithm?
- Counting operations
- Big-O
- Complexity classes

How fast is this code?

```
int search(int x, int[] A, int n)
16
17
      for (int i = 0; i < n; i++)</pre>
18
19
         if (A[i] == x) {
20
           return i;
21
22
23
      return -1;
24
```

How fast is this code?



Need a better way to measure

Permanent

Independent of hardware or other running processes, etc.

General

- Applicable to a large class of programs/algorithms/problems
- Resources: time, space

Mathematically rigorous

Useful

- Not for actual run time,
- But help us select best algorithm for the task

How many statements are executed?

```
int search(int x, int[] A, int n)
16
       for (int i = 0; i < n; i++)</pre>
17
                                            If x is not in A...
18
19
         if (A[i] == x) {
                                            how times are these
20
           return i;
                                            statements executed?
21
                                                    i = 0
22
23
      return -1;
                                            1+1
                                                   i < n
24
                                                    if (A[i] == x)
                                                    <u>i</u>++
                                                    return -1
```

How many operations are executed?

```
int search(int x, int[] A, int n)
16
       for (int i = 0; i < n; i++)</pre>
                                            If x is not in A...
18
19
         if (A[i] == x) {
                                            how times operations are
20
           return i;
                                            executed?
21
                                               l = 0
22
                                              \gamma + 1 = 1 < n
23
       return -1;
24
                                                    if (A[i] == x)
                                                    <u>i</u>++
                                                            (=(+)
```

How many operations are executed?

How many program operations are required to compute:

- L2 norm of vector
- Vector dot product
- Frobenius norm of matrix
- Matrix-vector multiplication
- Matrix-matrix multiplication

```
def norm(a):
    ss=0
    for i in range(len(a)):
        ss = ss + a[i]*a[i]
    norm = np.sqrt(ss)
    return norm
```

Operations:

- Arithmetic operations (e.g. + or **)
- Logical operations (e.g., and)
- Comparison operations (e.g., <=)
- Structure accessing operations (e.g. array indexing like A[i])
- Simple assignment such as copying a value into a variable
- Calls to library functions that don't depend on size of input (e.g., print)
- Control Statements (e.g. if X>5)

Be careful with function calls that scale with the size of the input