# Linear Regression / Optimization for ML

Matt Gormley
Lecture 7
Feb. 6, 2019

# Q&A
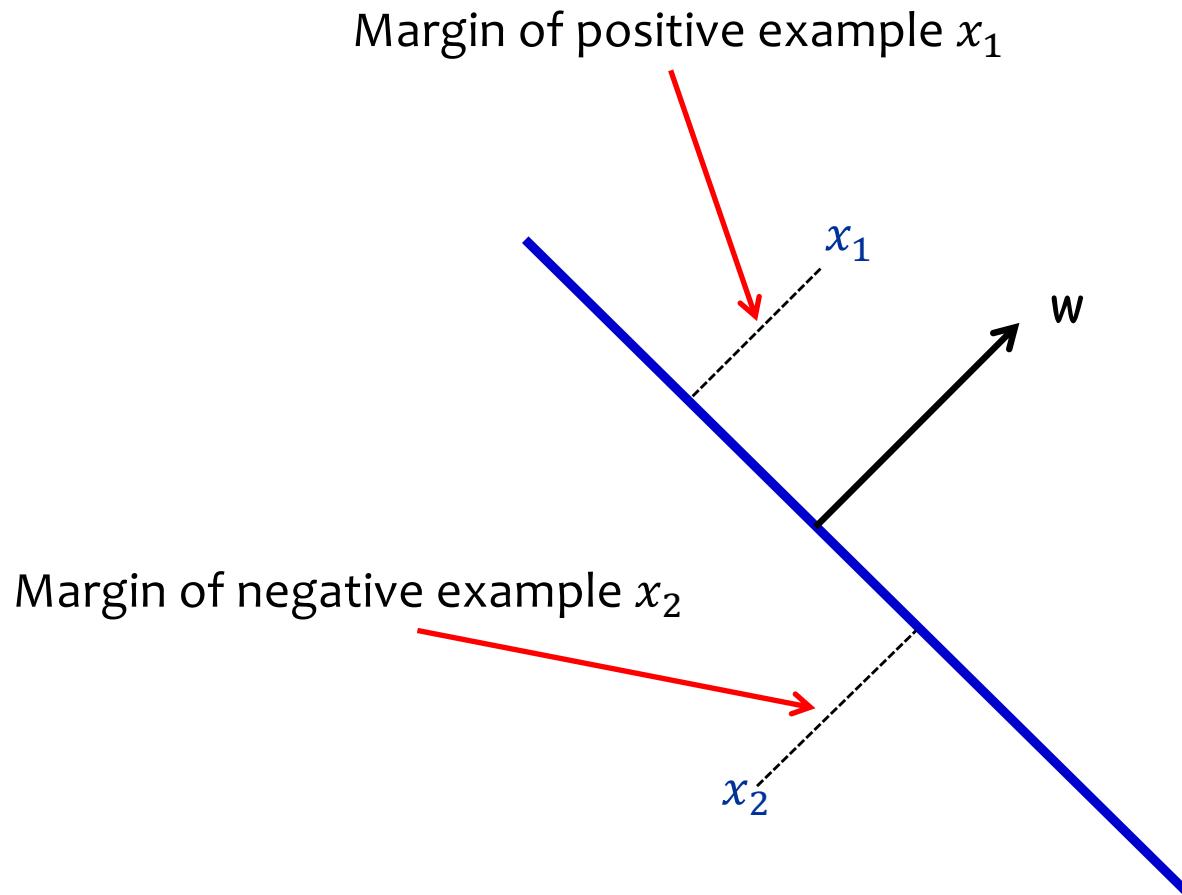
# Reminders

- **Homework 2: Decision Trees**
  - **Out: Wed, Jan 23**
  - **Due: Wed, Feb 6 at 11:59pm**

- **Homework 3: KNN, Perceptron, Lin.Reg.**
  - **Out: Wed, Feb 6**
  - **Due: Fri, Feb 15 at 11:59pm**

- **Today's In-Class Poll**
  - **http://p7.mlcourse.org**

# ANALYSIS OF PERCEPTRON

# Geometric Margin

**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)

Margin of positive example $x_1$

$x_1$

w
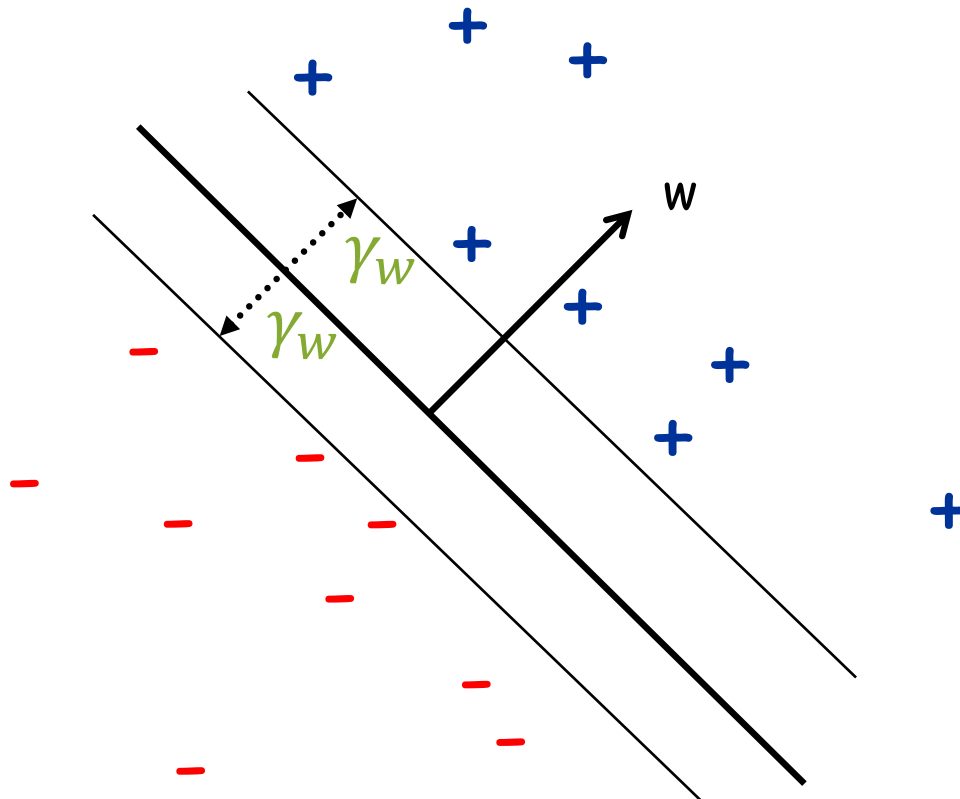
Margin of negative example $x_2$

$x_2$

# Geometric Margin

**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 $\gamma_w$ of a set of examples $S$ wrt a linear separator $w$ is the smallest margin over points $x \in S$.
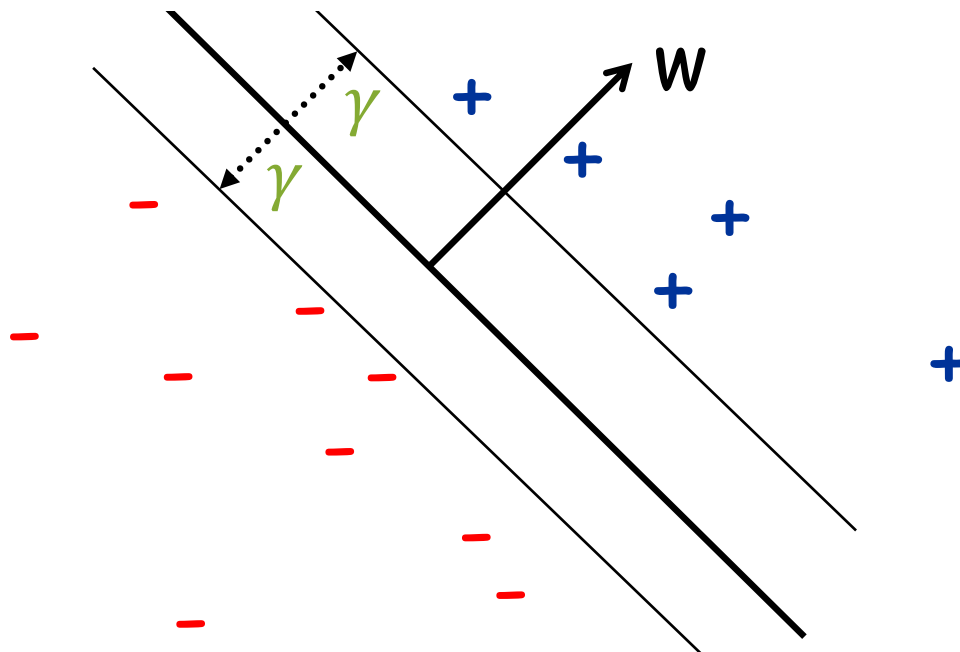
# Geometric Margin

**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 $\gamma_w$ of a set of examples $S$ wrt a linear separator $w$ is the smallest margin over points $x \in S$.
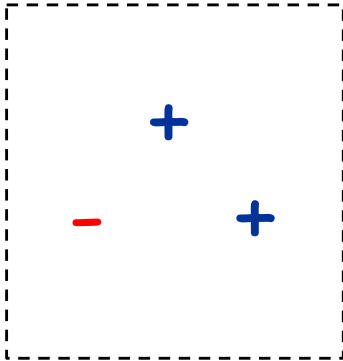
**Definition:** The margin $\gamma$ of a set of examples $S$ is the maximum $\gamma_w$ over all linear separators $w$.
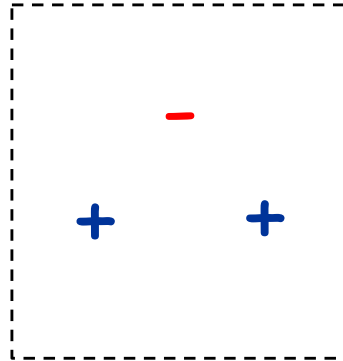
# Linear Separability

*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

Case 1:

$+$

$-$ $+$

Case 2:

$-$

$+$ $+$

Case 3:

$+$

$+$ $+$

Case 4:

$+$ $-$

$-$ $+$

# Analysis: Perceptron

## Perceptron Mistake Bound

**Guarantee:** If data has margin $\gamma$ 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.)
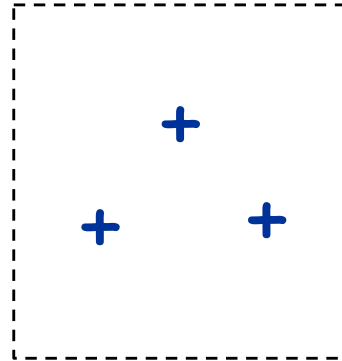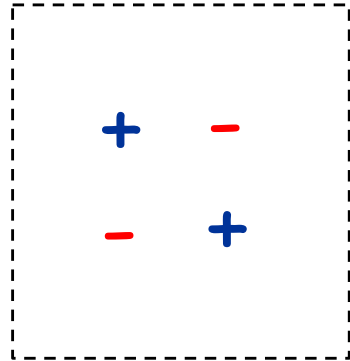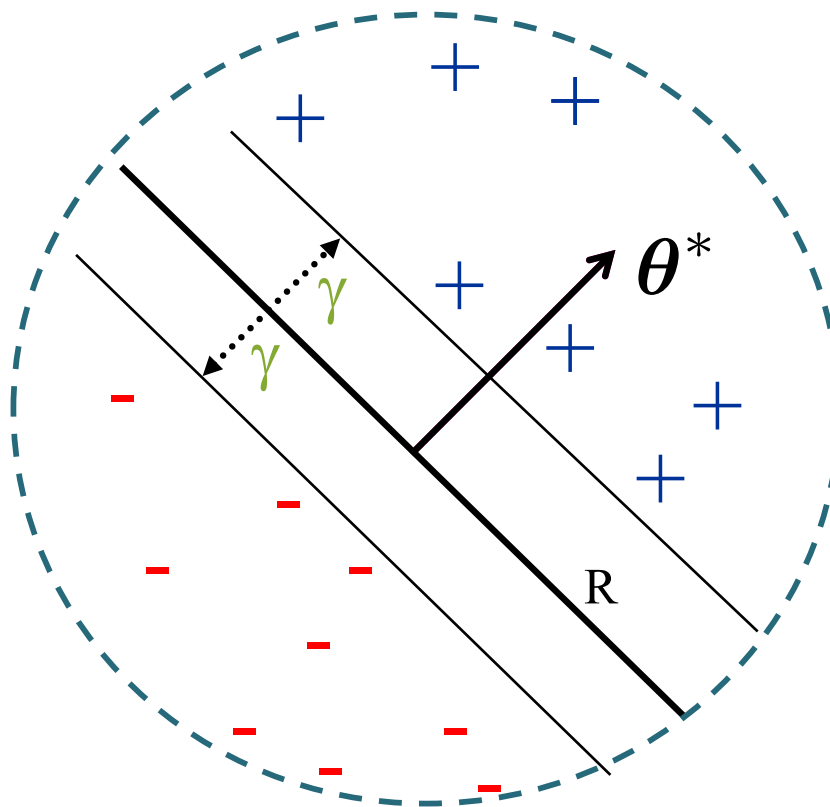
# Analysis: Perceptron

**Perceptron Mistake Bound**

**Guarantee:** If data has margin $\gamma$ 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 (batch) 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.

# Analysis: Perceptron

**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:* $||x^{(i)}|| \leq R$

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

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

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

Figure from Nina Balcan

# Analysis: Percep[tron]

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

**Perceptron Mistake Boun[d]**
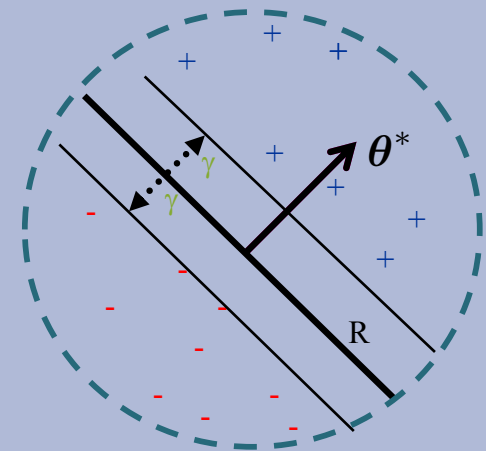
**Theorem 0.1** (Block (1962), Novikoff (19[...])

*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:* $\exists \boldsymbol{\theta}^* \text{ s.t. } ||\boldsymbol{\theta}^*|| = 1$ *and*
   $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

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

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

Figure from Nina Balcan

# Analysis: Perceptron

**Proof of Perceptron Mistake Bound:**

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

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

# Analysis: Perceptron

**Theorem 0.1** (Block (1962), Novikoff (1962)).
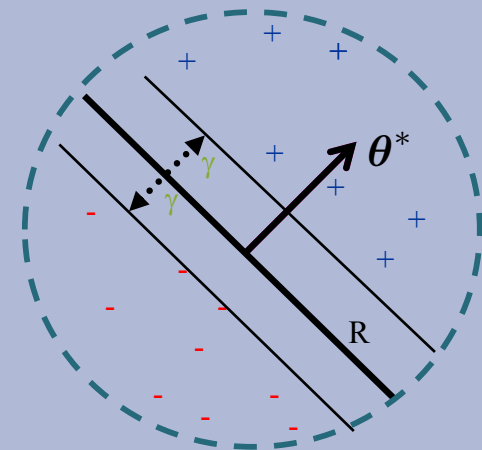*Given dataset:* $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$.
*Suppose:*

    1. *Finite size inputs:* $||x^{(i)}|| \leq R$

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

*Then: The number of mistakes made by the Perceptron algorithm on this dataset is*
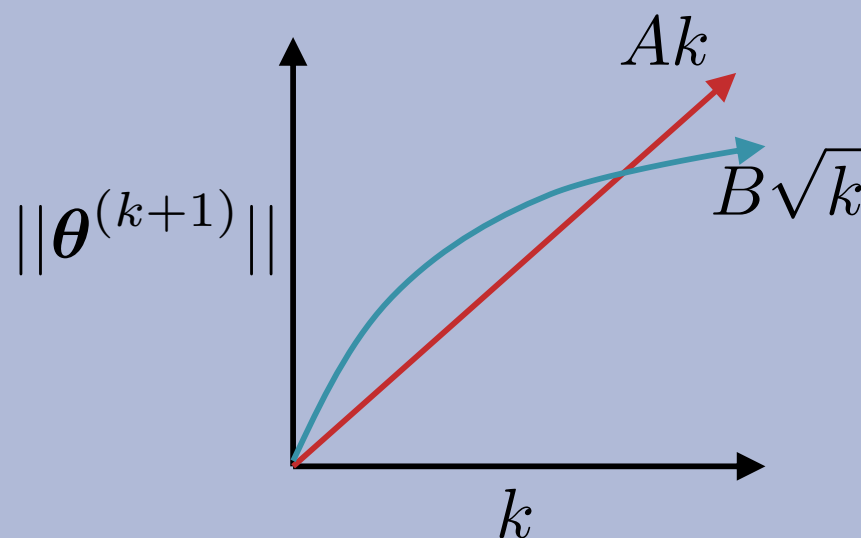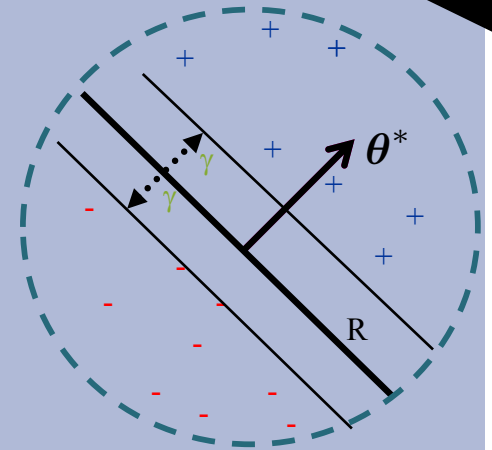
$$k \leq (R/\gamma)^2$$

---

**Algorithm 1** Perceptron Learning Algorithm (Online)

1:   **procedure** PERCEPTRON($\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \ldots\}$)
2:       $\boldsymbol{\theta} \leftarrow \mathbf{0}$, $k = 1$         ▷ Initialize parameters
3:       **for** $i \in \{1, 2, \ldots\}$ **do**         ▷ For each example
4:         **if** $y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$ **then**         ▷ If mistake
5:           $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)}$         ▷ Update parameters
6:         $k \leftarrow k + 1$
7:       **return** $\boldsymbol{\theta}$

15

# 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

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

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

17

# Analysis: Perceptron

**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 $k$th 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 $(y^{(i)})^2 = 1$

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

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

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

# Analysis: Perceptron

**Proof of Perceptron Mistake Bound:**

Part 3: Combining the bounds finishes the proof.

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

$$\Rightarrow k \leq (R/\gamma)^2$$

The total number of mistakes must be less than this

# Analysis: Perceptron

**What if the data is *not* linearly separable?**

1.  Perceptron will **not converge** in this case (it can't!)
2.  However, Freund & Schapire (1999) show that by projecting the points (hypothetically) into a higher dimensional space, we can achieve a similar bound on the number of mistakes made on **one pass** through the sequence of examples

**Theorem 2.** *Let $\langle(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m)\rangle$ be a sequence of labeled examples with $\|\mathbf{x}_i\| \leq R$. Let $\mathbf{u}$ be any vector with $\|\mathbf{u}\| = 1$ and let $\gamma > 0$. Define the deviation of each example as*

$$d_i = \max\{0, \gamma - y_i(\mathbf{u} \cdot \mathbf{x}_i)\},$$

*and define $D = \sqrt{\sum_{i=1}^m d_i^2}$. Then the number of mistakes of the online perceptron algorithm on this sequence is bounded by*

$$\left(\frac{R + D}{\gamma}\right)^2.$$

# Perceptron Exercises

**Question:**

*Unlike Decision Trees and K-Nearest Neighbors, the Perceptron algorithm **does not suffer from overfitting** because it does not have any hyperparameters that could be over-tuned on the training data.*

A. *True*
B. *False*
C. *True and False*

# Summary: Perceptron

- Perceptron is a **linear classifier**

- **Simple learning algorithm**: when a mistake is made, add / subtract the features

- Perceptron will converge if the data are **linearly separable**, it will **not** converge if the data are **linearly inseparable**

- For linearly separable and inseparable data, we can **bound the number of mistakes** (geometric argument)

- **Extensions** support nonlinear separators and structured prediction

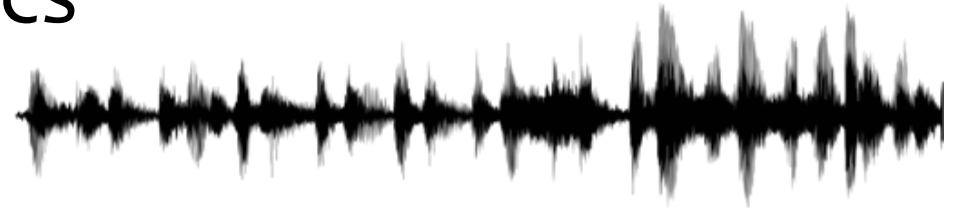# Perceptron Learning Objectives
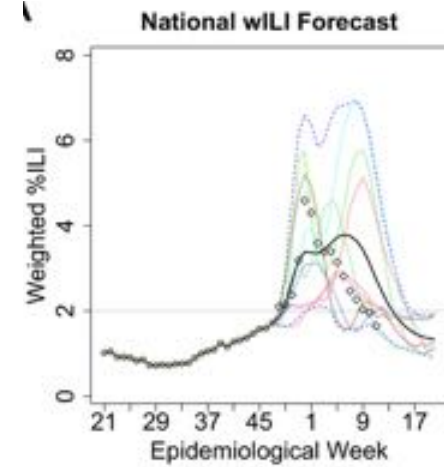
*You should be able to…*

- Explain the difference between online learning and batch learning
- Implement the perceptron algorithm for binary classification [CIML]
- Determine whether the perceptron algorithm will converge based on properties of the dataset, and the limitations of the convergence guarantees
- Describe the inductive bias of perceptron and the limitations of linear models
- Draw the decision boundary of a linear model
- Identify whether a dataset is linearly separable or not
- Defend the use of a bias term in perceptron

# LINEAR REGRESSION AS FUNCTION APPROXIMATION

# Regression



National wILI Forecast

Example Applications:

- Stock price prediction

- Forecasting epidemics

- Speech synthesis

- Generation of images (e.g. *Deep Dream*)

- Predicting the number of tourists on Machu Picchu on a given day

# Regression Problems

*Chalkboard*

- Definition of Regression

- Linear functions

- Residuals

- Notation trick: fold in the intercept

# Linear Regression as Function Approximation

*Chalkboard*

- Objective function: Mean squared error
- Hypothesis space: Linear Functions

# OPTIMIZATION IN CLOSED FORM

# Optimization for ML

Not quite the same setting as other fields...
- Function we are optimizing might not be the true goal
  (e.g. likelihood vs generalization error)
- Precision might not matter
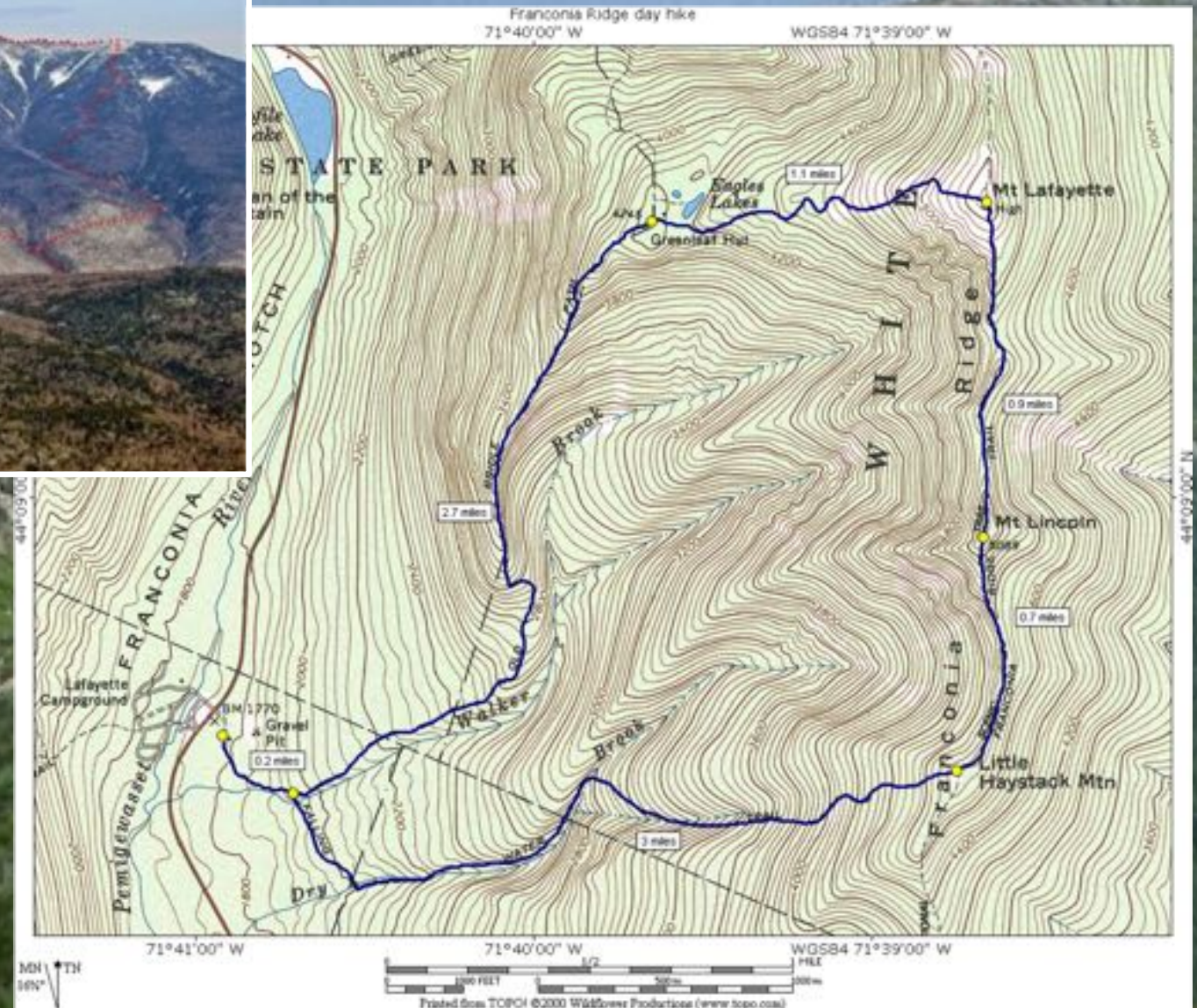  (e.g. data is noisy, so optimal up to 1e-16 might not help)
- Stopping early can help generalization error
  (i.e. "early stopping" is a technique for regularization – discussed more next time)

# Topographical Maps

# Topographical Maps
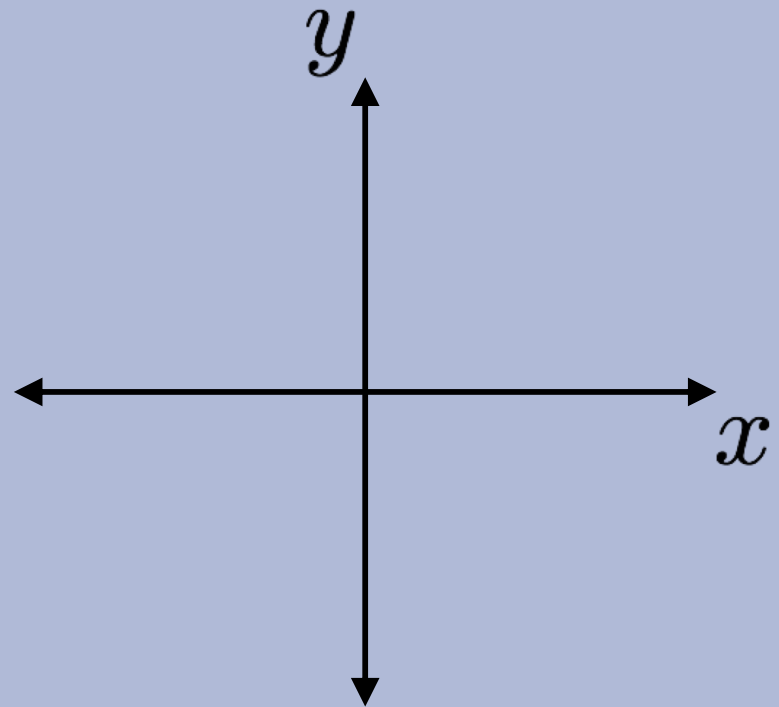
# Calculus and Optimization

**In-Class Exercise**

Plot three functions:

1. $f(x) = x^3 - x$

2. $f'(x) = \dfrac{\partial y}{\partial x}$

3. $f''(x) = \dfrac{\partial^2 y}{\partial x^2}$

**Answer Here:**

# Optimization for ML

*Chalkboard*
- Unconstrained optimization
- Convex, concave, nonconvex
- Derivatives
- Zero derivatives
- Gradient and Hessian

# Optimization: Closed form solutions

*Chalkboard*
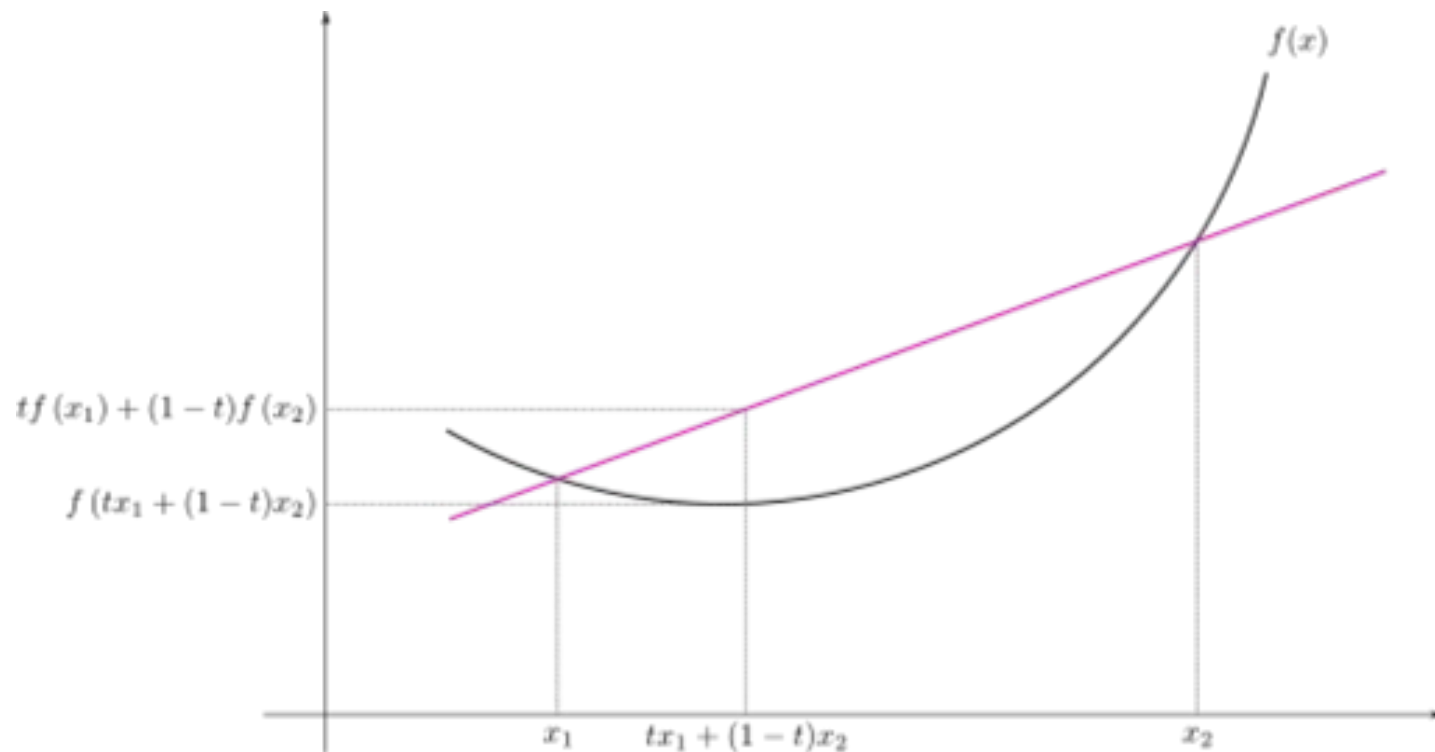
- Example: 1-D function
- Example: higher dimensions

# Convexity

Function $f : \mathbb{R}^M \to \mathbb{R}$ is **convex**
if $\forall\ \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \le t \le 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \le tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

There is only one local optimum if the function is *convex*

# Convexity

Function $f : \mathbb{R}^M \to \mathbb{R}$ is **convex**
if $\forall\ \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

There is only one local optimum if the function is *convex*

The **Mean Squared Error** function, which we will minimize for learning the parameters of Linear Regression, **is convex!**

41

# CLOSED FORM SOLUTION FOR LINEAR REGRESSION

# Optimization for Linear Regression

*Chalkboard*

- Closed-form (Normal Equations)
- Computational complexity of Closed-form Solution
- Stability of Closed-form Solution

# Function Approximation

*Chalkboard*
- The Big Picture