Machine Learning (CSE 446): Perceptron Convergence

C 2018

University of Washington cse446-staff@cs.washington.edu

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Review

Happy Medium?

Decision trees (that aren't too deep): use relatively few features to classify.

K-nearest neighbors: all features weighted equally.

Today: use all features, but weight them.

For today's lecture, assume that $y \in \{-1, +1\}$ instead of $\{0, 1\}$, and that $\mathbf{x} \in \mathbb{R}^d$.

Inspiration from Neurons

Image from Wikimedia Commons.



Input signals come in through dendrites, output signal passes out through the axon.

Perceptron Learning Algorithm

g= sign (w.x +b) **Data**: $D = \langle (\mathbf{x}_n, y_n) \rangle_{n=1}^N$, number of epochs E **Result**: weights \mathbf{w} and bias binitialize: $\mathbf{w} = \mathbf{0}$ and $\mathbf{b} = 0$: for $e \in \{1, ..., E\}$ do for $n \in \{1, \ldots, N\}$, in random order do # predict f g = -1, y = 2 $\hat{y} = \operatorname{sign}\left(\mathbf{w} \cdot \mathbf{x}_n + \mathbf{b}\right);$ if $\hat{y} \neq y_n$ then いん ジャジ # update $\begin{vmatrix} \mathbf{w} \leftarrow \mathbf{w} + y_n \\ b \leftarrow b + y_n; \end{vmatrix}$ if g=1, y=-1 end end end return w, b

Algorithm 1: PERCEPTRONTRAIN

Linear Decision Boundary



Linear Decision Boundary



Interpretation of Weight Values

What does it mean when

- ▶ $w_1 = 100?$
- ▶ $w_2 = -1?$
- $\blacktriangleright w_3 = 0?$

What if ||w|| is "large"? Sensitivity issues Syou night think large Hull is more 'complicated "

Today

What would we like to do?

• Optimization problem: find a classifier which minimizes the classification loss.

Some W

- ► The perceptron algorithm can be viewed as trying to do this...
- ▶ Problem: (in general) this is an NP-Hard problem.
- Let's still try to understand it...

This is the general approach of loss function minimization: find parameters which make our training error ¹small' (and which also generalizes)

training

When does the perceptron not converge?



Linear Separability

A dataset $D = \langle (\mathbf{x}_n, y_n) \rangle_{n=1}^N$ is **linearly separable** if there exists some linear classifier (defined by \mathbf{w}, b) such that, for all $n, y_n = \text{sign}(\mathbf{w}, \mathbf{x}_n + b)$.

 $\|\chi\| \leq 1$

If data are separable, (without loss of generality) can scale so that:



Perceptron Convergence

Due to Rosenblatt (1958).

Theorem: Suppose data are scaled so that $\|\mathbf{x}_i\|_2 \leq 1$. Assume D is linearly separable, and let be \mathbf{w}_* be a separator with "margin 1" Then the perceptron algorithm will converge in at most $\|\mathbf{w}_*\|^2$ epochs.

- Let \mathbf{w}_t be the param at "iteration" t; $\mathbf{w}_0 = 0$
- "A Mistake Lemma": At iteration t

If we make a mistake, $\|\mathbf{w}_{t+1} - \mathbf{w}_*\|^2 = \|\mathbf{w}_t - \mathbf{w}_*\|^2$ If we do make a mistake, $\|\mathbf{w}_{t+1} - \mathbf{w}_*\|^2 \le \|\mathbf{w}_t - \mathbf{w}_*\|^2 - 1$

► The theorem directly follows from this lemma. Why?

 $\hat{y} = -1, \ y = 1$ Proof of the "Mistake Lemma" $= \chi \cdot \chi$ $W_{\pm\pm1} = W_{\pm} \neq \chi$ $= \|w_{t+x} - w^{*}\|^{2} = \|w_{t} - w^{*} + \chi\|^{2}$ $\|W_{LL} - W^*\|^2$ $= ||w_{t} - w_{x}||^{2} + 2(w_{t} - w^{*}) \cdot X + ||X||^{2}$ $= ||w_t - w_x||^2 + 2(w_t - w^x) \cdot \chi + 1$ $\leq ||w_{t} - w_{*}/|^{2} + 2(-1) + 1$ hext paye $\leq \|W_{\xi} - w_{\star}\|^2 - ($

<ロ> < (回) < (回) < (目) < (目) < (目) と (目) (1) /

Proof of the "Mistake Lemma" (more scratch space)

<ロト < 回 > < 三 > < 三 > < 三 > 三 の < ぐ 11/13

- Suppose w¹, w⁴, w¹⁰, w¹¹ ... are the parameters right after we updated (e.g. after we made a mistake).
- Idea: instead of using the final w^t to classify, we classify with a majority vote using w¹, w⁴, w¹⁰, w¹¹...
- ► Why?

See CIML for details: Implementation and variants.

Let $\mathbf{w}^{(e,n)}$ and $b^{(e,n)}$ be the parameters after updating based on the *n*th example on epoch *e*.

$$\hat{y} = \operatorname{sign}\left(\sum_{e=1}^{E}\sum_{n=1}^{N}\operatorname{sign}(\mathbf{w}^{(e,n)} \cdot \mathbf{x} + b^{(e,n)})\right)$$



<ロ> < (日)、< (日)、< (目)、< (目)、< (目)、< (日)、< (日)、</td>13/13



(ロ)・(回)・(三)・(三)・(三)・(○)へ(○) 13/13



<ロ> < (日)、< (日)、< (目)、< (目)、< (目)、< (日)、< (日)、</td>13/13





<ロ > < 部 > < 言 > < 言 > こ > うへで 13/13



<ロ> < (回)、<(回)、<(目)、<(目)、<(目)、<(目)、(日)、<(13/13)</td>

Frank Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–408, 1958.