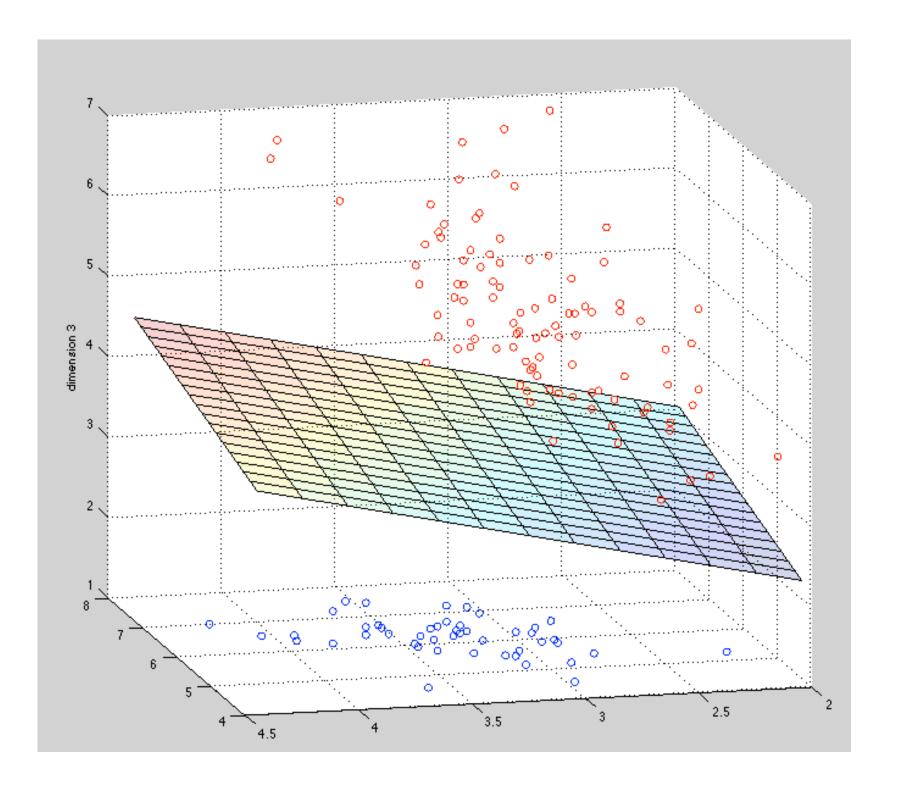
Perceptron

Machine Learning 10-601B
Seyoung Kim

Many of these slides are derived from William Cohen. Thanks!

Perceptron

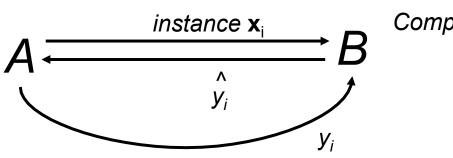
- Logistic regression is a linear classifier
- Another famous linear classifier
 - The perceptron



Probabilistic vs Margin-based Learning

- It's not all probabilities: many other types of analysis are used
- We also want to
 - capture geometric intuitions about what makes learning hard or easy
 - analyze performance worst-case settings
- This particular analysis is simple enough to give some insight into "margin" learning
- See: Freund & Schapire, 1998

[Rosenblatt, 1957]

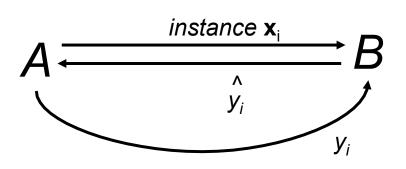


Compute: $\hat{y}_i = \text{sign}(\mathbf{v}_k \cdot \mathbf{x}_i)$ $\mathbf{v}_k \cdot \mathbf{x}_i$

If mistake: $\mathbf{v}_{k+1} = \mathbf{v}_k + y_i \mathbf{x}_i$

- On-line setting: data samples arrive one sample at a time
 - Adversary A provides student B with an instance x
 - Student B predicts a class (+1, -1) according to a simple linear classifier: sign($\mathbf{v}_k \cdot \mathbf{x}$)
 - Adversary gives student the answer (+1,-1) for that instance
- •Will do a *worst-case* analysis of the mistakes made by the student over *any* sequence of instances from the adversary
 - ... that follow a few rules

[Rosenblatt, 1957]

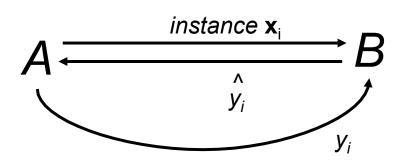


Compute:
$$\hat{y}_i = \text{sign}(\mathbf{v}_k \cdot \mathbf{x}_i)$$

If mistake:
$$\mathbf{v}_{k+1} = \mathbf{v}_k + y_i \mathbf{x}_i$$

- Amazingly simple algorithm
- Quite effective
- Very easy to *understand* if you do a little linear algebra

[Rosenblatt, 1957]



Compute: $\hat{y}_i = \text{sign}(\mathbf{v}_k \cdot \mathbf{x}_i)$

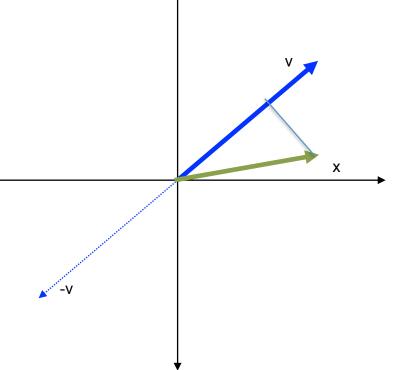
If mistake: $\mathbf{v}_{k+1} = \mathbf{v}_k + y_i \mathbf{x}_i$

• Recall dot product definition:

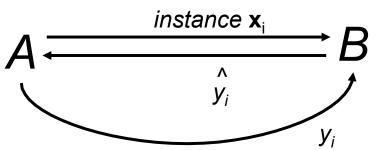
$$\mathbf{x} \bullet \mathbf{v} = \sum_{\mathbf{i}} x_i v_i$$

- •and intuition:
 - project vector x onto vector v
 - dot product is the distance from the origin to that projection

So why does this algorithm make sense?



[Rosenblatt, 1957]



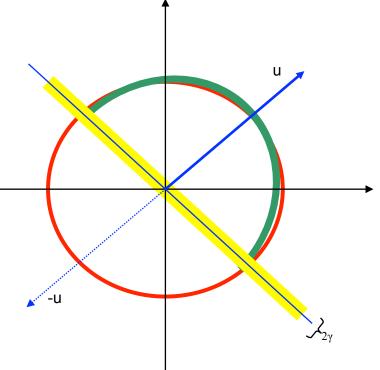
Compute: $\hat{y}_i = \text{sign}(\mathbf{v}_k \cdot \mathbf{x}_i)$

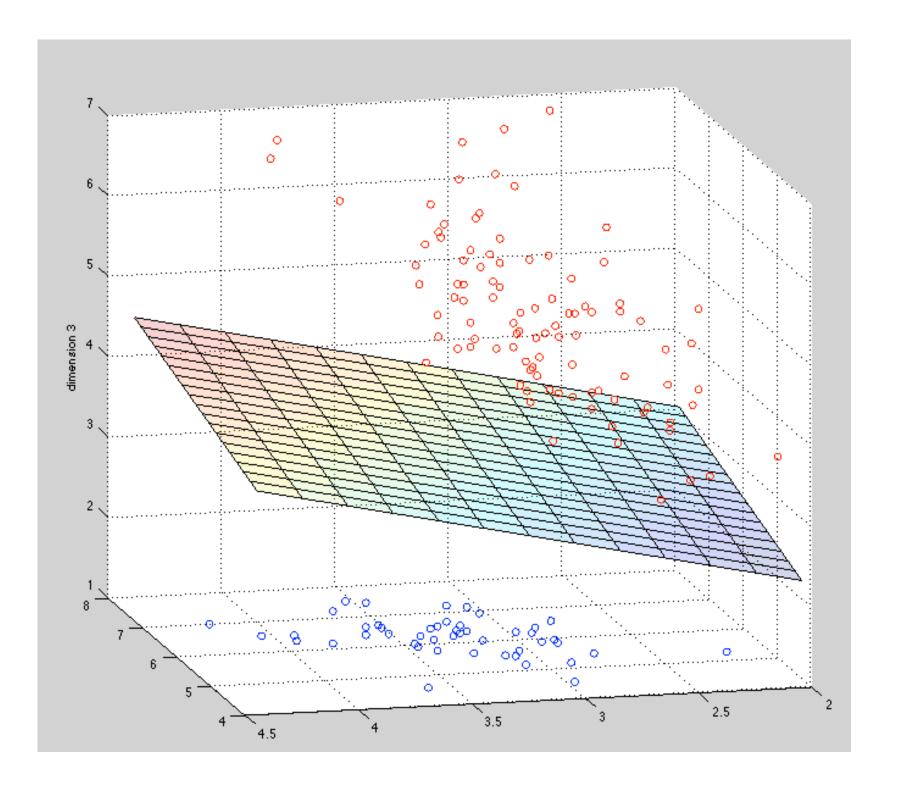
If mistake: $\mathbf{v}_{k+1} = \mathbf{v}_k + y_i \mathbf{x}_i$

- Amazingly simple algorithm
- Quite effective
- Very easy to understand if you do a little linear algebra

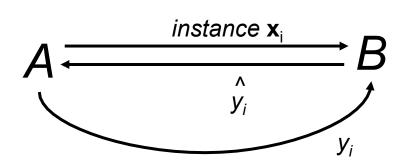


- Examples are not too "big"
- There is a "good" answer -- i.e. a line that clearly separates the pos/neg examples





[Rosenblatt, 1957]



Compute: $y_i = sign(\mathbf{v}_k \cdot \mathbf{x}_i)$

If mistake: $\mathbf{v}_{k+1} = \mathbf{v}_k + y_i \mathbf{x}_i$

Rule 1: Radius R: A must provide

examples "near the origin"

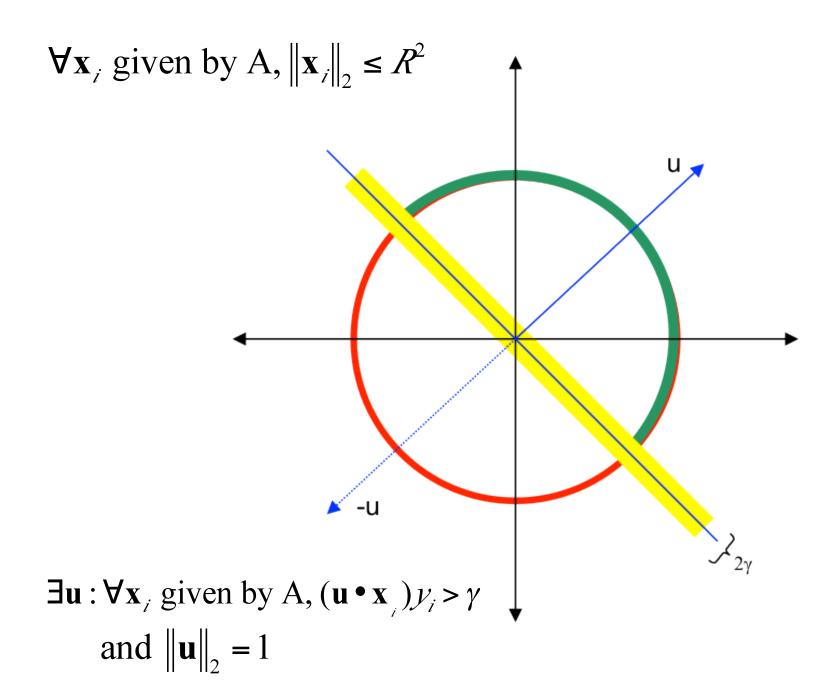
Rule 2: Margin
$$\gamma$$
: A must provide examples that can be separated with some vector \mathbf{u} with margin $\gamma>0$ and unit norm

$$\forall \mathbf{x}_{i} \text{ given by A, } \|\mathbf{x}_{i}\|_{2}^{2} \leq R^{2}$$

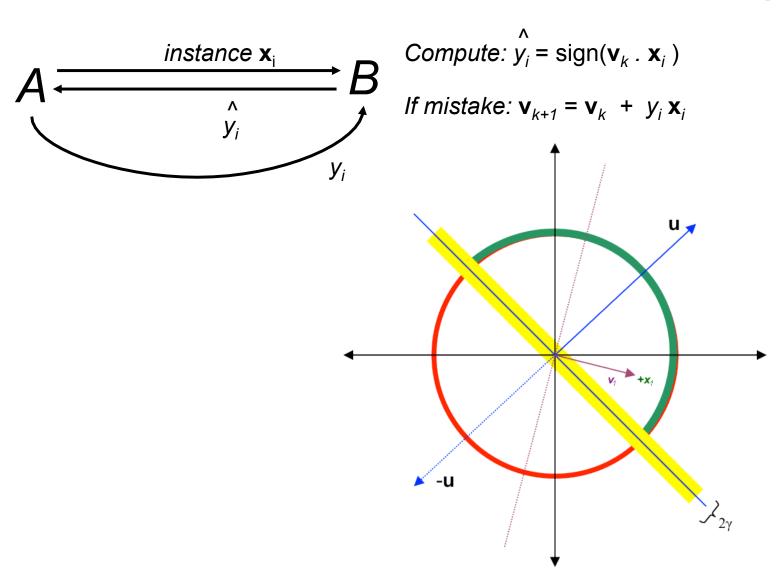
$$\|\mathbf{x}\|_{2} = \sqrt{(x_{1}^{2} + ... + x_{n}^{2})}$$

$$\exists \mathbf{u} : \forall \mathbf{x}_i \text{ given by A, } (\mathbf{u} \cdot \mathbf{x}_i) y_i > \gamma$$

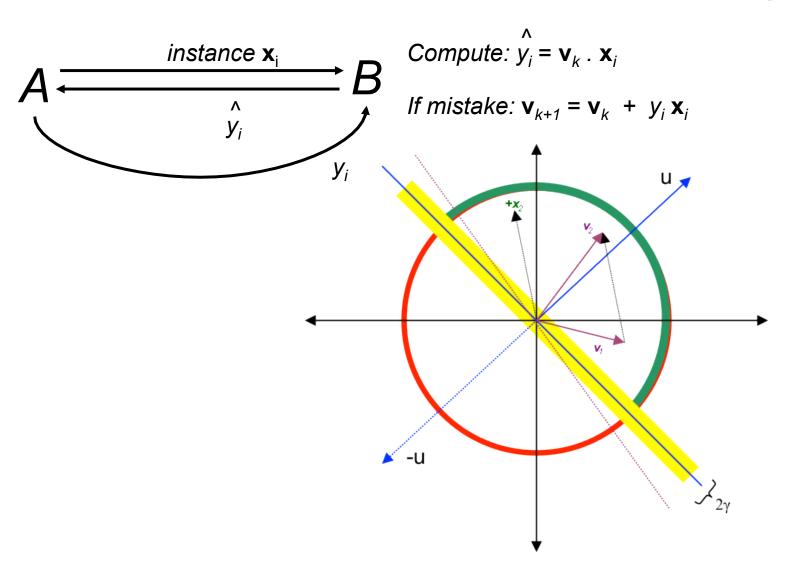
and $\|\mathbf{u}\|_2 = 1$



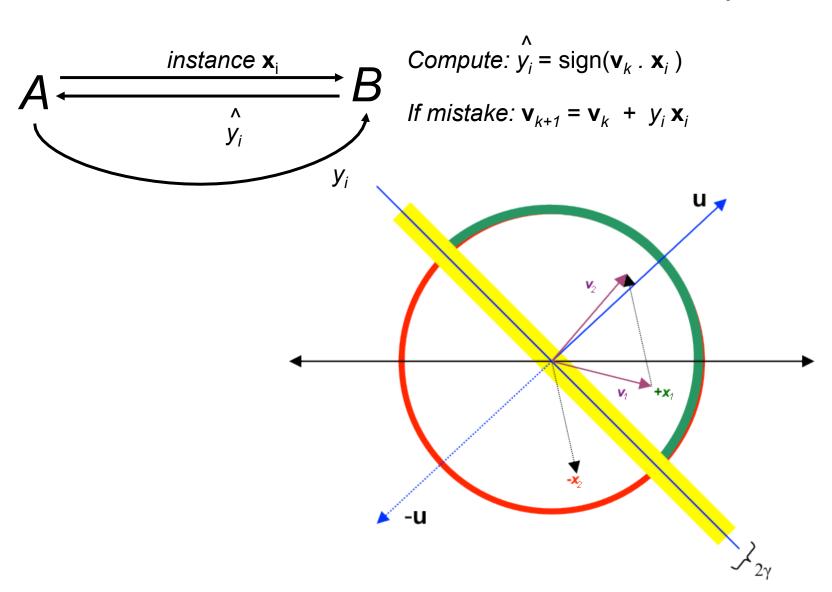
The perceptron: after one positive x_i



The perceptron: after two positive x_i

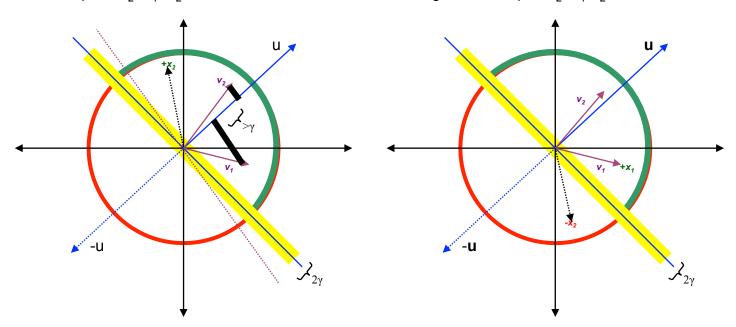


The perceptron: after one pos + one neg x_i



The guess $\mathbf{v_2}$ after the two positive examples: $\mathbf{v_2} = \mathbf{v_1} + \mathbf{x_2}$

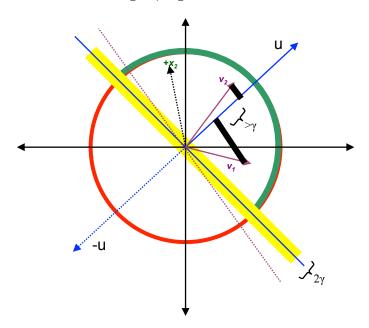
The guess $\mathbf{v_2}$ after the one positive and one negative example: $\mathbf{v_2} = \mathbf{v_1} - \mathbf{x_2}$



Lemma 1: the dot product between \mathbf{v}_k and \mathbf{u} increases with each mistake by at least γ : i.e.,

$$\forall k : \mathbf{v}_k \cdot \mathbf{u} \ge k \gamma$$

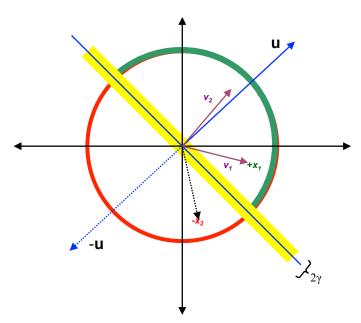
The guess $\mathbf{v_2}$ after the two positive examples: $\mathbf{v_2} = \mathbf{v_1} + \mathbf{x_2}$



Lemma 1: the dot product between \mathbf{v}_k and \mathbf{u} increases with each mistake by at least γ : i.e.,

$$\forall k : \mathbf{v}_k \cdot \mathbf{u} \ge k \gamma$$

The guess $\mathbf{v_2}$ after the one positive and one negative example: $\mathbf{v_2} = \mathbf{v_1} - \mathbf{x_2}$



$$\mathbf{v}_{k+1} \cdot \mathbf{u} = (\mathbf{v}_k + y_i \mathbf{x}_i) \cdot \mathbf{u}$$

$$\mathbf{v}_{k+1} \cdot \mathbf{u} = (\mathbf{v}_k \cdot \mathbf{u}) + y_i (\mathbf{x}_i \cdot \mathbf{u})$$

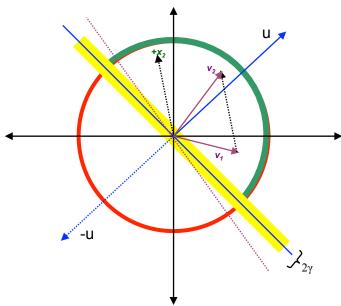
$$\mathbf{v}_{k+1} \cdot \mathbf{u} \ge (\mathbf{v}_k \cdot \mathbf{u}) + \gamma$$

$$\mathbf{v}_k \cdot \mathbf{u} \ge k\gamma$$

$$\mathbf{v}_k \cdot \mathbf{u} \ge k\gamma$$

$$\mathbf{v}_k \cdot \mathbf{u} \ge k\gamma$$

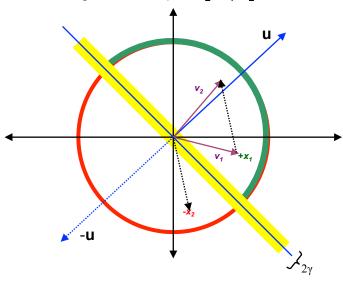
(3a) The guess \mathbf{v}_2 after the two positive examples: $\mathbf{v}_2 = \mathbf{v}_1 + \mathbf{x}_2$



Lemma 2: The norm of \mathbf{v}_k grows slowly with each mistake, i.e.,

$$\forall k, \left\| \mathbf{v}_{k} \right\|_{2}^{2} \le kR^{2}$$

(3b) The guess $\mathbf{v_2}$ after the one positive and one negative example: $\mathbf{v_2} = \mathbf{v_1} - \mathbf{x_2}$



$$\begin{aligned} \mathbf{v}_{k+1} \cdot \mathbf{v}_{k+1} &= (\mathbf{v}_k + y_i \mathbf{x}_i) \cdot (\mathbf{v}_k + y_i \mathbf{x}_i) \\ \left\| \mathbf{v}_{k+1} \right\|_2^2 &= \left\| \mathbf{v}_k \right\|_2^2 + 2y_i \mathbf{v}_k \cdot \mathbf{x}_i + y_i^2 \left\| \mathbf{x}_i \right\|_2^2 \\ \left\| \mathbf{v}_{k+1} \right\|_2^2 &\leq \left\| \mathbf{v}_k \right\|_2^2 + 1 \left\| \mathbf{x}_i \right\|_2^2 \\ \left\| \mathbf{v}_{k+1} \right\|_2^2 &\leq \left\| \mathbf{v}_k \right\|_2^2 + R^2 \end{aligned}$$
 Always negative,

 $\forall \mathbf{x}_{i} \text{ given by A}, \|\mathbf{x}_{i}\|_{2}^{2} \leq R^{2}$ SO ... $\|\mathbf{v}_{k}\|_{2}^{2} \leq kR^{2}$

Always negative since it was a mistake

Lemma 1: the dot product between \mathbf{v}_{k} and \mathbf{u} increases with each mistake by at last y: i.e.,

Lemma 2: The norm of \mathbf{v}_k grows slowly with each mistake, i.e.,

$$\forall k: \mathbf{v}_{k} \cdot \mathbf{u} \ge k\gamma \qquad \forall k, \|\mathbf{v}_{k}\|$$

$$k\gamma \le \mathbf{v}_{k} \cdot \mathbf{u} \quad \text{and} \quad \|\mathbf{v}_{k}\|_{2}^{2} \le kR^{2}$$

$$k^{2}\gamma^{2} \le \|\mathbf{v}_{k} \cdot \mathbf{u}\|_{2}^{2} \quad \text{and} \quad \|\mathbf{v}_{k}\|_{2}^{2} \le kR^{2}$$

$$k^{2}\gamma^{2} \le \|\mathbf{v}_{k}\|_{2}^{2} \cdot \|\mathbf{u}\|_{2}^{2} \quad \text{and} \quad \|\mathbf{v}_{k}\|_{2}^{2} \le kR^{2}$$

$$k^{2}\gamma^{2} \le \|\mathbf{v}_{k}\|_{2}^{2} \quad \text{and} \quad \|\mathbf{v}_{k}\|_{2}^{2} \le kR^{2}$$

$$k^{2}\gamma^{2} \le \|\mathbf{v}_{k}\|_{2}^{2} \le kR^{2}$$

$$k^{2}\gamma^{2} \le kR^{2}$$

$$k^{2}\gamma^{2} \le kR^{2}$$

$$k < \left(\frac{R}{\gamma}\right)^{2}$$

$$\forall k, \left\| \mathbf{v}_{k} \right\|_{2}^{2} \le kR^{2}$$

...and
$$\|\mathbf{u}\|_{2} = 1$$

Summary

- We have shown that
 - If: exists a **u** with unit norm that has margin γ on examples in the seq $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots$
 - *Then*: the perceptron algorithm makes $< R^2/\gamma^2$ mistakes on the sequence (where R $>= ||\mathbf{x}_i||$)
 - Independent of dimension of the data (!)
- We don't know what happens if the data's not separable



On-line to batch learning

Imagine we run the on-line perceptron and see this result.

_			<u> </u>	
i	guess	input	result	
1	\mathbf{v}_0	\mathbf{x}_1	X (a mistake)	Which v _i should we use?
2	${f v}_1$	\mathbf{x}_2	$\sqrt{\text{(correct!)}}$	·
3	\mathbf{v}_1	\mathbf{x}_3	\checkmark	Maybe the <i>last</i> one?
4	${f v}_1$	\mathbf{x}_4	X (a mistake)	Here it's never gotten any
5	\mathbf{v}_2	X_5	\checkmark	test cases right!
6	\mathbf{v}_2	\mathbf{x}_6	\checkmark	(Experimentally, the classifiers move around a lot.)
7	\mathbf{v}_2	\mathbf{x}_7	\checkmark	Maybe the "best one"?
8	\mathbf{v}_2	\mathbf{x}_8	X	•
9	\mathbf{v}_3	\mathbf{x}_9	\checkmark	But we "improved" it with
10	\mathbf{v}_3	\mathbf{x}_{10}	X	later mistakes

$$P(\text{error in } \mathbf{x}) = \sum_{k} P(\text{error on } \mathbf{x}|\text{picked } \mathbf{v}_{k}) P(\text{picked } \mathbf{v}_{k})$$

$$= \sum_{k} \frac{1}{m_{k}} \frac{m_{k}}{m} = \sum_{k} \frac{1}{m} = \frac{k}{m}$$

Imagine we run the on-line perceptron and see this result.

i	guess	input	result
1	\mathbf{v}_0	\mathbf{x}_1	X (a mistake)
2	${f v}_1$	\mathbf{x}_2	$\sqrt{\text{(correct!)}}$
3	\mathbf{v}_1	\mathbf{x}_3	\checkmark
4	\mathbf{v}_1	\mathbf{x}_4	X (a mistake)
5	\mathbf{v}_2	\mathbf{X}_5	\checkmark
6	\mathbf{v}_2	\mathbf{x}_6	\checkmark
7	\mathbf{v}_2	\mathbf{x}_7	\checkmark
8	\mathbf{v}_2	\mathbf{x}_8	X
9	\mathbf{v}_3	\mathbf{x}_9	\checkmark
10	\mathbf{v}_3	\mathbf{x}_{10}	X

- Pick a v_k at random according to m_k/m, the fraction of examples it was used for.
- 2. Predict using the \mathbf{v}_k you just picked.
- 3. Better: use a deterministic approximation to this: a sum of the \mathbf{v}_k 's, weighted by m_k/m

From Freund & Schapire, 1998: Classifying digits with VP

