

## Features and kernels

### Online Perceptron

**input** Labeled examples  $\{(x_i, y_i)\}_{i=1}^n$  from  $\mathbb{R}^d \times \{-1, +1\}$ .

```

1: initialize  $\hat{w}_1 := \mathbf{0}$ .
2: for  $t = 1, 2, \dots, n$  do
3:   if  $y_t \langle \hat{w}_t, x_t \rangle \leq 0$  then
4:      $\hat{w}_{t+1} := \hat{w}_t + y_t x_t$ .
5:   else
6:      $\hat{w}_{t+1} := \hat{w}_t$ 
7:   end if
8: end for
9: return  $\hat{w}_{n+1}$ .

```

**Theorem:** If  $R := \max_{t \in \{1, \dots, n\}} \|x_t\|_2$ , and  $w_* \in \mathbb{R}^d$  satisfies

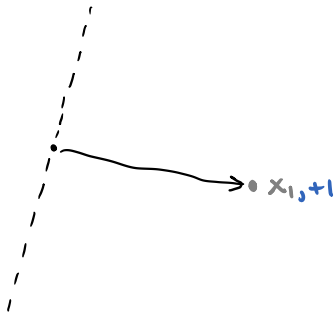
$$y \langle w_*, x_t \rangle \geq 1 \quad \text{for all } (x_t, y_t),$$

then Online Perceptron makes at most  $\|w_*\|_2^2 \cdot R^2$  mistakes (and updates).

1 / 29

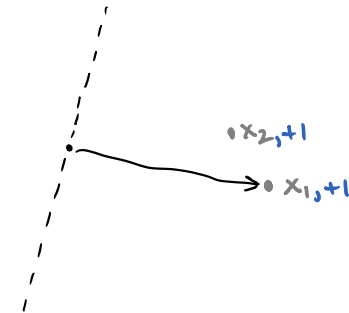
2 / 29

## Example run of Online Perceptron



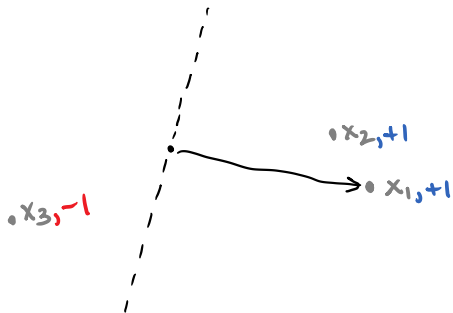
3 / 29

## Example run of Online Perceptron



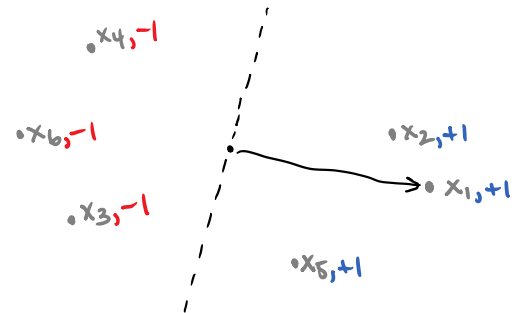
4 / 29

Example run of Online Perceptron



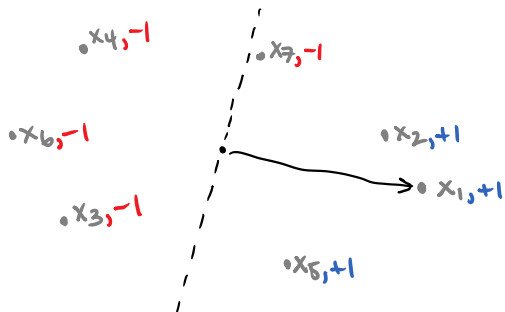
5 / 29

Example run of Online Perceptron



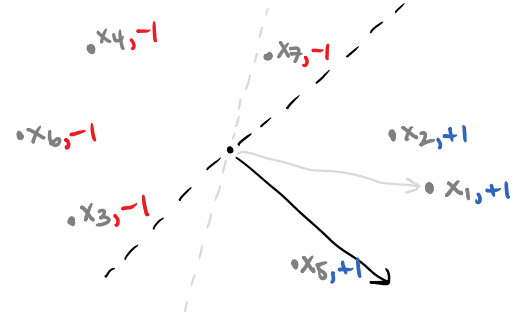
6 / 29

Example run of Online Perceptron



7 / 29

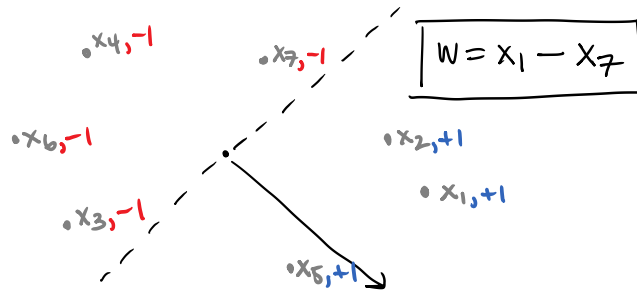
Example run of Online Perceptron



8 / 29

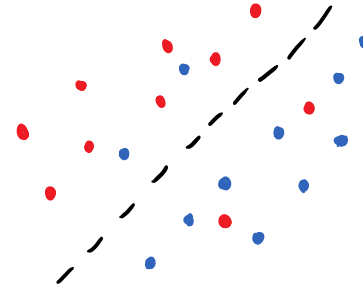
## What to do if data is not linearly separable

### Example run of Online Perceptron

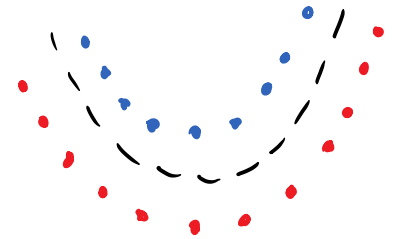


9 / 29

Noise: even Bayes classifier not perfect



Bayes classifier not (approx.) linear

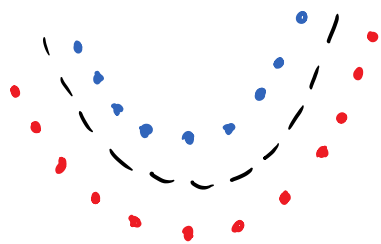


10 / 29

## Adding new features

Original feature vector:  $\mathbf{x} = (x_1, x_2)$ .

New feature vector:  $\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$ .



Decision boundary is non-linear in  $\mathbf{x}$ , but is linear in  $\phi(\mathbf{x})$ .

11 / 29

## Getting the most out of linear classifiers

Often, with a “good” set of features, a linear classifier can well-approximate the Bayes classifier.

**Two approaches:**

1. Think very hard and carefully about which features to use.
2. Use all features that come to mind.

12 / 29

## The kitchen sink of features

**Example:** document classification

► **Word features:**

$\mathbb{1}\{\text{"aardvark" appears}\}, \mathbb{1}\{\text{"abacus" appears}\}, \dots, \mathbb{1}\{\text{"zygote" appears}\}$

► **Bi-gram features:**

$\mathbb{1}\{\text{"bank deposit" appears}\}, \mathbb{1}\{\text{"river bank" appears}\}, \dots$

► **Tri-gram features:**

$\mathbb{1}\{\text{"New York City" appears}\}, \mathbb{1}\{\text{"wherefore art thou" appears}\}, \dots$

**Example:** new features from old features  $\mathbf{x} \in \mathbb{R}^d$

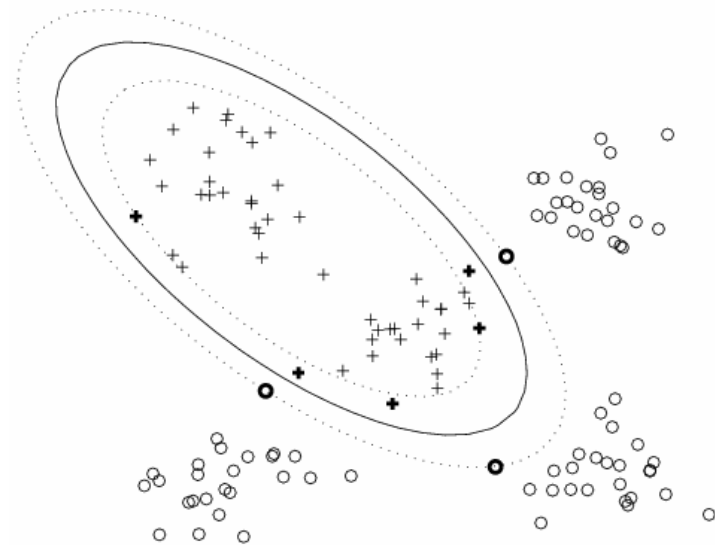
► **Pairwise interactions:**

$$(x_1x_2, x_1x_3, \dots, x_1x_d, x_2x_3, \dots, x_{d-1}x_d) \in \mathbb{R}^{\binom{d}{2}}$$

► Etc.

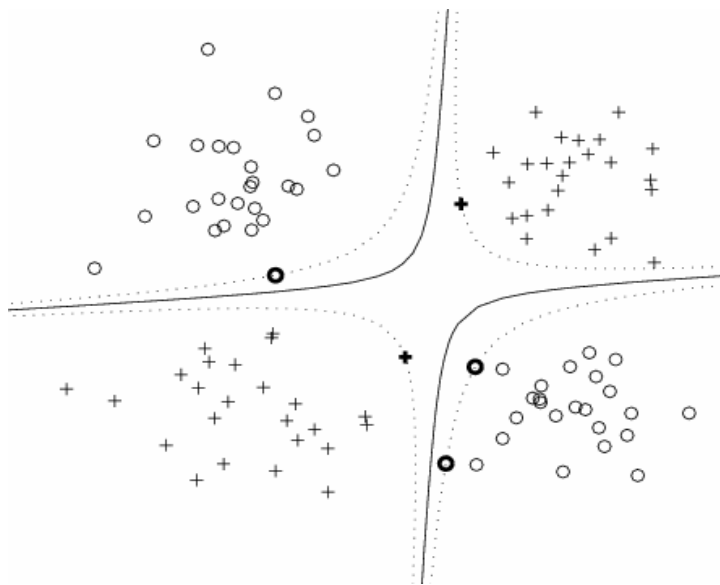
13 / 29

## All degree $\leq 2$ interaction features



14 / 29

## All degree $\leq 2$ interaction features



15 / 29

## Learning with the kitchen sink of features

- Let  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$  be the mapping to the expanded feature space (e.g., all original features and all deg. 2 interactions, so  $D = \Omega(d^2)$ ).
- Learn linear classifier  $f_w: \mathbb{R}^D \rightarrow \{\pm 1\}$  (i.e., learn a weight vector  $w \in \mathbb{R}^D$ ) using data with expanded features  $\{(\phi(\mathbf{x}_i), y_i)\}_{i=1}^n$ .
- **Caveat:** can be computationally expensive to do this directly if  $D$  is large. Naïvely: takes  $\Omega(D)$  time to even make a prediction.

16 / 29

- **Recall:** Online Perceptron weight vector (using expanded features) is

$$\mathbf{w} = \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x})$$

where  $\mathcal{M}$  is the subset of labeled examples where Online Perceptron made a mistake during training.

- **Prediction** using Online Perceptron weight vector on new point  $\mathbf{z} \in \mathbb{R}^d$ :

$$\langle \mathbf{w}, \phi(\mathbf{z}) \rangle = \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle.$$

**Computational cost:**  $|\mathcal{M}| \times$  time to compute inner product  $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$ .  
Sometimes this can be faster than  $O(D)$ .

- $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{1+2d+\binom{d}{2}}$ , where

$$\phi(\mathbf{x}) = \left( 1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_1x_d, \dots, \sqrt{2}x_{d-1}x_d \right)$$

(Don't mind the  $\sqrt{2}$ 's.)

- **Computing  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$  in  $O(d)$  time:**

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = (1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^2.$$

- Much better than  $\Omega(d^2)$ .

- $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{2^d}$ , where

$$\phi(\mathbf{x}) = \left( \prod_{i \in S} x_i : S \subseteq \{1, 2, \dots, d\} \right)$$

- **Computing  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$  in  $O(d)$  time:**

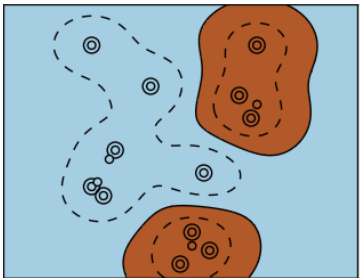
$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \prod_{i=1}^d (1 + x_i x'_i).$$

- Much better than  $\Omega(2^d)$ .

For any  $\sigma > 0$ , there is an infinite feature expansion  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^\infty$  (involving Hermite polynomials of all orders) such that

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}\right),$$

which can be computed in  $O(d)$  time.



(This is called the **Gaussian kernel** with bandwidth  $\sigma$ .)

A **kernel function**  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a symmetric function with the following property:

*For any  $x_1, x_2, \dots, x_n \in \mathcal{X}$ , the  $n \times n$  matrix whose  $(i, j)$ -th entry is  $K(x_i, x_j)$  is positive semidefinite.*

For any kernel  $K$ , there is a feature mapping  $\phi: \mathcal{X} \rightarrow \mathbb{H}$  such that

$$\langle \phi(x), \phi(x') \rangle = K(x, x'),$$

( $\mathbb{H}$  is a Hilbert space—i.e., a special kind of inner product space—called the *Reproducing Kernel Hilbert Space corresponding to  $K$* .)

►  $\phi: \text{Strings} \rightarrow \mathbb{N}^{\text{Strings}}$ , where

$$\phi(x) = (\text{number of times } s \text{ appears in } x : s \in \text{Strings})$$

$$K(x, x') = \langle \phi(x), \phi(x') \rangle = \text{measure of similarity between strings.}$$

► **Computing  $K(x, x')$ :**

For each substring  $s$  of  $x$ , count how many times  $s$  appears in  $x'$  and add to total.

Dynamic programming in time  $O(\text{length}(x) \times \text{length}(x'))$ .

► **Implicit representation:**

$$w = \sum_{(x,y) \in \mathcal{M}} y \phi(x)$$

- Never explicitly form weight vector  $w$ .
- Instead, store all labeled examples in  $\mathcal{M}$ .
- Whenever need to compute  $\langle w, \phi(z) \rangle$  for new point  $z$ , iterate over examples in  $\mathcal{M}$  to compute

$$\langle w, \phi(z) \rangle = \sum_{(x,y) \in \mathcal{M}} y \langle \phi(x), \phi(z) \rangle = \sum_{(x,y) \in \mathcal{M}} y K(x, z).$$

- Focus on designing good kernels (rather than feature maps), which means designing good **similarity functions**.
- Lots of ways to construct kernels.  
(E.g., combine existing kernels.)
- Lots of algorithms can be / have been “kernelized” (like Perceptron).

# Experimental results on OCR

- ▶ OCR digits data, binary classification problem: distinguish “9” from other digits.
- ▶ # training examples: 60000 (about 6000 are from class “9”).
- ▶ Test error rates using Kernelized Averaged Perceptron (similar to Voted Perceptron).

# passes	0.1	1	2	3	4	10
Linear kernel	0.045	0.039	0.038	0.038	0.038	0.037
Degree 2	0.024	0.012	0.010	0.010	0.009	0.009
Degree 4	0.020	0.009	0.008	0.007	0.007	0.006

25 / 29

# More computational issues

- ▶ **Implicit representation:**

$$w = \sum_{(x,y) \in \mathcal{M}} y \phi(x)$$

- ▶ Number of mistakes  $|\mathcal{M}|$  could be  $\Omega(n)$ .
- ▶ Computing predictions as expensive as brute-force NN search.
- ▶ Training algorithms quite slow (e.g.,  $\Omega(n^2)$ ).

26 / 29

# Kernel approximations

## Many ways to try to speed-up kernel methods using approximations.

- ▶ Limit number of examples used to represent weight vector.  
“Nystrom approximation”  
“Budgeted Perceptron”

- ▶ Use *explicit* feature maps  $z: \mathbb{R}^d \rightarrow \mathbb{R}^m$  such that

$$\langle z(x), z(x') \rangle \approx K(x, x').$$

“Random projections / feature hashing”  
“Random kitchen sinks”

27 / 29

# Experimental results on e-mail data

- ▶ Spam data set (4601 e-mail messages, 39.4% are spam).
- ▶  $\mathcal{Y} = \{\text{spam}, \text{not spam}\}$ ,  $\mathcal{X} = \mathbb{R}^{57}$  (message features)
- ▶ # training examples: 3065, # test examples: 1536
- ▶ Test error rates
  - ▶ **Decision tree learning:** 9.3%
  - ▶ **Averaged Perceptron** (128 passes): 8.27%
  - ▶ **Random Kitchen Sink Averaged Perceptron** (64 passes): 6.12% (approximates Gaussian kernel)

28 / 29

## Key takeaways

1. Linear classifiers only as good as given feature representation.
2. Explicit feature expansion (sometimes okay!)
3. Kernel trick: sometimes never need  $\phi(\mathbf{x})$  directly, but only via  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ , computed quickly as  $K(\mathbf{x}, \mathbf{x}')$ .
4. Kernel approach: “similarity engineering” rather than “feature engineering”
5. High-level idea of using kernel approximations.