# Winnowing Subspaces

#### Manfred Warmuth

University of California - Santa Cruz

ICML 2007, June 21 - 2007

### Plots by Dima Kuzmin

-∢∃>

# To Winnow



MW ()

3

### How I do it



MW ()

#### To winnow:

- to remove (chaff from grain) by a current of air
- to get rid of (something undesirable or unwanted)
- From Anglo-Saxon "windwian"

[L]

- Online alg. for learning disjunctions
- Mistake bound logarithmic in number of features winnows a large number of features
- Will morph into algorithm for winnowing subspaces

# Disjunctions as linear threshold functions

- 2 out of 5 literal monotone disjunction  $v_1 \lor v_3$
- Represented as  $\mathbf{d} = (1, 0, 1, 0, 0)^{\top}$
- Label for instance  $\mathbf{x} = (0, 1, 1, 0, 0)^{\top}$

$$\begin{cases} +1 & \text{if } \mathbf{d} \cdot \mathbf{x} \geq \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

• Alg. receives sequence of examples online

$$(\mathbf{x}_1, y_1) \ (\mathbf{x}_2, y_2), \ \dots, \ (\mathbf{x}_T, y_T)$$

instances  $[0,1]^n$ , labels and predictions are  $\pm 1$ 

Initialize 
$$\mathbf{w}_1 = w_0 \ (1, 1, \dots 1)^{\top}$$
  
for  $t = 1$  to  $\mathcal{T}$  do  
Receive instance  $\mathbf{x}_t \in [0, 1]^n$   
Predict with  
 $\hat{y}_t = \begin{cases} +1 & \text{if } \mathbf{w}_t \cdot \mathbf{x}_t \ge \theta \\ -1 & \text{otherwise} \end{cases}$   
Receive label  $y_t \in \{+1, -1\}$   
Update  
 $w_{t+1,i} = \begin{cases} w_{t,i} & \text{if no mistake} \\ w_{t,i} \ e^{\eta y_t \times_{t,i}} & \text{if mistake} \end{cases}$ 

### end for

Like perceptron alg., except multiplicative update

伺下 イヨト イヨト

- Thm If examples consistent with k out of n literal monotone disjunction, then properly tuned Winnow make at most O(k ln n) mistakes
- Mistake bound logarithmic in dimension *n*
- Perceptron alg. can make  $\Omega(n k)$  mistakes

- What variables?
- What dot product?
- What corresponds to disjunctions?
- What happens to the exponential form of weights soft max?



伺下 イヨト イヨト

# Visualization of symmetric matrices

• As ellipses - affine transformations of the unit ball



• Ellipse = {Wu :  $||u||_2 = 1$ }

### Ellipses cont.



• Eigenvectors form the axes and eigenvalues their lengths

-





Degenerate ellipses

- One eigenvalue one
- All others zero

4

### Linear combinations and mixtures of dyads

### • Symmetric matrices are linear combinations of dyads



• Positive definite matrices

Eigenvalues are non-negative

• Density matrices are mixtures of dyads Eigenvalues form probability vector (

• Many mixtures lead to same density matrix

$$0.2 \longrightarrow + 0.3 + 0.5 = \begin{pmatrix} 0.35 & 0.15 \\ 0.15 & 0.65 \end{pmatrix} = = 0.29 + 0.71$$

- There always exists a decomposition into *n* dyads that correspond to eigenvectors
- Uncertainty about dyad expressed as density matrix
- We have a Bayes rule for density matrices

### Variance

 View the symmetric positive definite matrix W as a covariance matrix of some random cost vector c ∈ R<sup>n</sup>

$$\mathbf{W} = \mathbb{E}\left( (\mathbf{c} - \mathbb{E}(\mathbf{c}))(\mathbf{c} - \mathbb{E}(\mathbf{c}))^{\top} 
ight)$$

 $\bullet\,$  The variance along any vector u is

$$\mathbf{V}(\mathbf{c}^{\top}\mathbf{u}) = \mathbb{E}\left(\left(\mathbf{c}^{\top}\mathbf{u} - \mathbb{E}(\mathbf{c}^{\top}\mathbf{u})\right)^{2}\right)$$
$$= \mathbf{u}^{\top}\underbrace{\mathbb{E}\left((\mathbf{c} - \mathbb{E}(\mathbf{c}))(\mathbf{c} - \mathbb{E}(\mathbf{c}))^{\top}\right)}_{\mathbf{W}}\mathbf{u}$$

Variance as trace

$$\mathbf{u}^{\top}\mathbf{W}\mathbf{u} = \operatorname{tr}(\mathbf{u}^{\top}\mathbf{W}\mathbf{u}) = \operatorname{tr}(\mathbf{W}\ \mathbf{u}\mathbf{u}^{\top}) \geq 0$$

### Plotting the variance

(u<sup>T</sup>Wu)u Curve of the ellipse is plot of vector  $\mathbf{W}\mathbf{u}$ , where  $\mathbf{u}$  is unit vector The outer figure eight is direction **u** times the variance  $\mathbf{u}^{\mathsf{T}}\mathbf{W}\mathbf{u}$ For an eigenvector, this variance equals the eigenvalue and touches the ellipse

MW ()

# 3 dimensional variance plots



# What dot product?

$$\operatorname{tr}(\mathbf{W} \mathbf{X}) = \operatorname{tr}(\sum_{i} \omega_{i} \mathbf{w}_{i} \mathbf{w}_{i}^{\top} \mathbf{X})$$
$$= \sum_{i} \omega_{i} \operatorname{tr}(\mathbf{w}_{i} \mathbf{w}_{i}^{\top} \mathbf{X})$$
$$= \sum_{i} \omega_{i} \underbrace{\mathbf{w}_{i}^{\top} \mathbf{X} \mathbf{w}_{i}}_{\text{variance along eigendirs}}$$

Measurement in quantum physics

- Dyad  $\mathbf{u}\mathbf{u}^{\top}$  is state
- Density matrix W is mixture state
- Instance matrix X is instrument
- $tr(\mathbf{W} \mathbf{X})$  is expected outcome



-

- 4 聞 と 4 直 と 4 直 と

# What corresponds to disjunctions

Disjunctions

$$(1,0,1,0,0)^{\top} \cdot (x_1,x_2,x_3,x_4,x_5)^{\top} = x_1 + x_3$$

Sum k components of **x** 

• Projections matrices  $\mathbf{P} = \sum_{i=1}^{k} \mathbf{u}_i \mathbf{u}_i^{\top}$ 

$$\operatorname{tr}(\mathsf{PX}) = \sum_{i=1}^{k} \mathsf{u}_{i}^{\top} \mathsf{X} \mathsf{u}_{i}$$

Sum variance along k directions

$$(\mathbf{X}_{1}, y_{1}), \ (\mathbf{X}_{2}, y_{2}), \ \ldots, \ (\mathbf{X}_{T}, y_{T}), \ \hat{y}_{T}$$

Label  $y_t$  is +1 if trace is at least  $\frac{1}{2}$  and -1 otherwise

MW ()

# Thresholding the trace for $\mathbf{x}\mathbf{x}^{\top}$ instances





MW ()

Winnowing Subspaces

22 / 24

# Symmetric Matrix Winnow

Initialize  $\mathbf{W}_1 = w_0 \prod_{n \times n}^{t}$ for t = 1 to T do Receive instance  $\underset{n \times n}{\mathbf{X}_t}$  with eigenvalues in [0, 1]Predict with  $\hat{y}_t = \begin{cases} +1 & \text{if } \operatorname{tr}(\mathbf{W}_t \mathbf{X}_t) \ge \theta \\ -1 & \text{otherwise} \end{cases}$ 

Receive label  $y_t$ Update

 $\mathbf{W}_{t+1} = \begin{cases} \mathbf{W}_t & \text{if no mistake} \\ \exp(\log \mathbf{W}_t + \eta \ y_t \mathbf{X}_t) & \text{if mistake} \end{cases}$ 

#### end for

**exp** and **log** are spectral functions In normalized version, trace normalized to one

MW ()

◆□▶ ◆□▶ ◆目▶ ◆目▶ 三回 のへの

- Same  $O(k \ln n)$  mistake bound if examples consistent with *k*-dimensional subspace
  - Dubbed free matrix lunch
- Generalizes to arbitrary matrices use mixtures of  $\boldsymbol{u}\boldsymbol{v}^{\top}$  and SVD
- Key tool in analysis is quantum relative entropy

[TRW]

 $\Delta(\mathbf{W},\mathbf{V}) = \operatorname{tr}(\mathbf{W}(\log \mathbf{W} - \log \mathbf{V}))$