

# RVMs for classification

Lecture 25

Consider  $K=2: t \in \{0, 1\}$ .

$$y(\vec{x}, \vec{w}) = \sigma(\vec{w}^T \cdot \vec{\psi}(\vec{x}))$$

We will use a separate  $\alpha_i$  for each  $w_i$ .

Given  $\vec{\alpha} = \underline{\alpha_1 \dots \alpha_M}$ , we have:

$$\log p(\vec{w} | \vec{t}, \vec{\alpha}) = \underbrace{\sum_{n=1}^N [t_n \log y_n + (1-t_n) \log(1-y_n)]}_{\log \mathcal{L}} -$$

$$\underbrace{\frac{1}{2} \vec{w}^T A \vec{w} + \text{const}(\vec{w})}_{\log(\text{prior})} \quad (*)$$

$$A = \begin{pmatrix} \alpha_1 & & 0 \\ & \ddots & \\ 0 & & \alpha_M \end{pmatrix}$$

Maximize (\*) w.r.t  $\vec{w}$  to get  $\vec{w}^*$ :

$$\frac{\partial}{\partial w_j} [t_n \log y_n + (1-t_n) \log(1-y_n)] =$$

$$= \underbrace{\left[ \frac{t_n}{y_n} \psi_j(\vec{x}_n) - \frac{(1-t_n)}{1-y_n} \psi_j(\vec{x}_n) \right]}_{y_n(1-y_n)} = \frac{t_n(1-y_n) - (1-t_n)y_n}{y_n(1-y_n)} \psi_j(\vec{x}_n) \quad (\ominus)$$

$$\uparrow \text{from } \sigma'(\vec{w}^T \cdot \vec{\psi}) \quad (\ominus) \psi_j(\vec{x}_n) (t_n - y_n) = \varphi_{jn}^T (t_n - y_n).$$

So,  $\nabla_{\vec{w}} \log p(\vec{w} | \vec{t}, \vec{\alpha}) \Big|_{\vec{w}^*} = \varphi^T (\vec{t} - \vec{y}) - A \vec{w} \Big|_{\vec{w}^*} = 0$ , so that

$$\vec{w}^* = A^{-1} \varphi^T (\vec{t} - \vec{y}).$$



Then  $\lambda_i (w_i^*)^2 = 1 - \lambda_i \sum_{ii}$ , or

$$\lambda_i^{\text{new}} = \frac{\gamma_i}{(w_i^*)^2}$$

update eq'n  
evaluated at  $\vec{\lambda}^{\text{old}}, \vec{w}^*$

Iterate to convergence between estimating  $\{\vec{w}^*, \vec{\lambda}\}$   $\Rightarrow$  obtain sparse solution, usually much sparser than SVMs

We can also analyze  $\lambda_i$  dependence more explicitly by using  $\log p(\vec{t} | \vec{\lambda}) \equiv L(\vec{\lambda})$  & splitting off the  $\lambda_i$ -dependent terms:  
 $L(\vec{\lambda}) = L(\vec{\lambda}_{-i}) + \lambda(\lambda_i)$  as before

Finally, RVMs generalize to  $K > 2$  w/out difficulties:

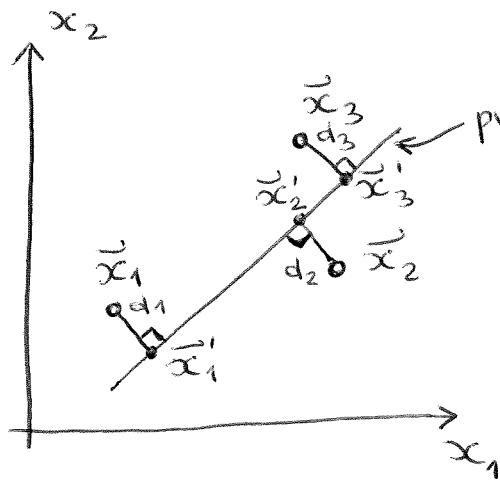
$$y_k(\vec{x}) = \frac{e^{a_k}}{\sum_{j=1}^K e^{a_j}}, \quad a_k = \vec{w}_k^T \vec{\phi}(\vec{x})$$

softmax f'n

# Principal component analysis (PCA)

Goals: dimensionality reduction, feature extraction, data visualization

$D=2$



principal subspace

either maximize the variance of  $\bar{x}'_1, \bar{x}'_2, \bar{x}'_3, \dots$  or minimize the  $\perp$  distances  $d_1, d_2, d_3, \dots$

## Max variance formulation

Consider  $\{\bar{x}_n\}_{n=1}^N$ , dimensionality  $D$   
 # components in  $\bar{x}_n$

Goal: ~~find~~ project  $\bar{x}_n$ 's onto a subspace with  $M < D$ , while maximizing the variance of the projected data.

Start with  $M=1$ : define  $\bar{u}_1$  s.t.  $\bar{u}_1^T \cdot \bar{u}_1 = 1$   
 $\underbrace{\quad}_{D \text{ dim's}}$

Then  $\bar{x}_n \Rightarrow x_n' = \bar{u}_1^T \cdot \bar{x}_n$ , and

$$\underbrace{\langle \bar{x}' \rangle}_{\text{mean}} = \frac{1}{N} \sum_{n=1}^N x_n' = \bar{u}_1^T \cdot \underbrace{\frac{1}{N} \sum_n \bar{x}_n}_{\langle \bar{x} \rangle} = \bar{u}_1^T \cdot \langle \bar{x} \rangle$$

$$\begin{aligned} \text{Var}\{x_n'\} &= \frac{1}{N} \sum_n (\bar{u}_1^T \cdot \bar{x}_n - \bar{u}_1^T \cdot \langle \bar{x} \rangle)^2 = \\ &= \frac{1}{N} \sum_n u_{1,k} [x_{n,k} x_{n,j} + \langle x \rangle_k \langle x \rangle_j - \\ &\quad - x_{n,k} \langle x \rangle_j - x_{n,j} \langle x \rangle_k] u_{1,j} \quad \textcircled{=} \end{aligned}$$

$$\textcircled{=} u_{1,k} \left\{ \frac{1}{N} \sum_n (x_{n,k} - \langle x \rangle_k)(x_{n,j} - \langle x \rangle_j) \right\} u_{1,j} =$$

" $S_{kj}$ "

$$= \vec{u}_1^T S \vec{u}_1$$

maximize this, subject to  $\vec{u}_1^T \vec{u}_1 = 1$ :

$$\mathcal{J} = \vec{u}_1^T S \vec{u}_1 + \lambda_1 (1 - \vec{u}_1^T \vec{u}_1)$$

↑ Lagrange multiplier

$$\frac{\partial \mathcal{J}}{\partial \vec{u}_1} = 2(S \vec{u}_1 - \lambda_1 \vec{u}_1) = 0 \Rightarrow S \vec{u}_1 = \underline{\underline{\lambda_1 \vec{u}_1}}$$

Thus  $\vec{u}_1$  is an eigenvector of  $S$ , with eigenvalue  $\lambda_1$ .

Note that  $\vec{u}_1^T S \vec{u}_1 = \lambda_1$ ,

the variance will be at max if  $\lambda_1$  is the largest eigenvalue of  $S$ ; corresponding  $\vec{u}_1$  is called the 1<sup>st</sup> principal component.

We can continue to project, s.t. for  $M$ -dim principal subspace the optimal projection is defined by  $M$  largest eigenvalues:

$$\lambda_1, \dots, \lambda_M$$

$$\vec{u}_1, \dots, \vec{u}_M \leftarrow \text{corresponding eigenvectors}$$

Proof by induction:

$M=1$  shown above.

For  $M > 1$ , assume this holds for  $M$  & prove it for  $M+1$ . Thus

$\vec{u}_1, \dots, \vec{u}_M$  are principal eigenvectors

and we are looking for  $\vec{u}_{M+1} \perp \vec{u}_1, \dots,$

$$\vec{u}_{M+1}^T \vec{u}_{M+1} = 1$$

$$\vec{u}_{M+1} \perp \vec{u}_M.$$

otherwise  $\vec{u}_{M+1}$  lies in the  $M$ -subspace

Since var in the  $\vec{u}_{M+1}$  direction is given by  $\vec{u}_{M+1}^T S \vec{u}_{M+1}$ , we need to maximize

$$\mathcal{J} = \vec{u}_{M+1}^T S \vec{u}_{M+1} + \lambda_{M+1} (1 - \vec{u}_{M+1}^T \vec{u}_{M+1}) + \sum_{i=1}^M \eta_i \vec{u}_{M+1}^T \vec{u}_i$$

$$\frac{\partial \mathcal{J}}{\partial \vec{u}_{M+1}} = 0 \Rightarrow 2S \vec{u}_{M+1} - 2\lambda_{M+1} \vec{u}_{M+1} + \sum_{i=1}^M \eta_i \vec{u}_i = 0.$$

$$2 \vec{u}_{M+1}^T S \vec{u}_{M+1} - 2\lambda_{M+1} \underbrace{\vec{u}_{M+1}^T \vec{u}_{M+1}}_1 + \sum_{i=1}^M \eta_i (\vec{u}_{M+1}^T \vec{u}_i) = 0 \quad \text{or}$$

$$\vec{u}_{M+1}^T S \vec{u}_{M+1} = \lambda_{M+1}$$

$\Downarrow$

$$S \vec{u}_{M+1} = \lambda_{M+1} \vec{u}_{M+1}$$

maximized when  $\lambda_{M+1}$  is the largest remaining eigenvalue

## Min-error formulation

Again, introduce  $\{\vec{u}_i\}_{i=1}^D$  s.t.

$$\vec{u}_j^T \cdot \vec{u}_i = \delta_{ij}$$

$$\text{Then } \vec{x}_n = \sum_{i=1}^D \alpha_{ni} \vec{u}_i = \sum_{i=1}^D (\vec{x}_n^T \cdot \vec{u}_i) \vec{u}_i$$

$$\text{Approximate } \vec{x}_n \Rightarrow \tilde{\vec{x}}_n = \sum_{i=1}^M z_{ni} \vec{u}_i + \sum_{i=M+1}^D b_i \vec{u}_i \quad \leftarrow \text{indep. of } n$$

Now, choose  $\{\vec{u}_i\}$ ,  $\{z_{ni}\}$ ,  $\{b_i\}$  to

$$\text{minimize } J = \frac{1}{N} \sum_{n=1}^N \|\vec{x}_n - \tilde{\vec{x}}_n\|^2.$$

$$\frac{\partial J}{\partial z_{nj}} = \frac{2}{N} \left( \vec{x}_n - \sum_{i=1}^M z_{ni} \vec{u}_i - \sum_{i=M+1}^D b_i \vec{u}_i \right)^T \cdot \vec{u}_j = 0 \quad 1 \leq j \leq M$$

$$= \frac{2}{N} (\vec{x}_n^T \cdot \vec{u}_j - z_{nj}) = 0 \Rightarrow z_{nj} = \vec{x}_n^T \cdot \vec{u}_j.$$

Likewise,

$$\frac{\partial J}{\partial b_i} = \frac{2}{N} \sum_{n=1}^N \left( \vec{x}_n - \sum_{i=1}^M z_{ni} \vec{u}_i - \sum_{i=M+1}^D b_i \vec{u}_i \right)^T \cdot \vec{u}_j = 0 \quad M+1 \leq j \leq D$$

$$= 2 \left( \frac{1}{N} \sum_{n=1}^N (\vec{x}_n^T \cdot \vec{u}_j) - b_j \right) = 0, \text{ or}$$

$$b_j = \langle \vec{x} \rangle^T \cdot \vec{u}_j.$$

Now,

$$\vec{x}_n - \tilde{\vec{x}}_n = \sum_{i=M+1}^D [(\vec{x}_n - \langle \vec{x} \rangle)^T \cdot \vec{u}_i] \vec{u}_i$$

"residual" lies in the subspace defined by  $\vec{u}_{M+1} \dots \vec{u}_D$

Furthermore,

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (\vec{x}_n^T \cdot \vec{u}_i - \langle \vec{x} \rangle^T \cdot \vec{u}_i)^2 =$$

$$= \sum_{i=M+1}^D \vec{u}_i^T S \vec{u}_i$$

— 0 —

Consider  $D=2, M=1$  case first:

minimize  $J = \vec{u}_2^T S \vec{u}_2$  subject to  $\vec{u}_2^T \cdot \vec{u}_2 = 1$ .

$$\tilde{J} = \vec{u}_2^T S \vec{u}_2 + \lambda_2 (1 - \vec{u}_2^T \vec{u}_2)$$

$$\frac{\partial \tilde{J}}{\partial \vec{u}_2} = 0 \Rightarrow \underline{S \vec{u}_2 = \lambda_2 \vec{u}_2} \Rightarrow J = \lambda_2$$

To minimize  $J$ , choose  $\lambda_2$  to be the smaller eigenvalue  $\Rightarrow$  the principal subspace is defined by  $\vec{u}_1$  corresponding to  $\lambda_1$ , the larger eigenvalue.



In general, for  $M < D$  solve

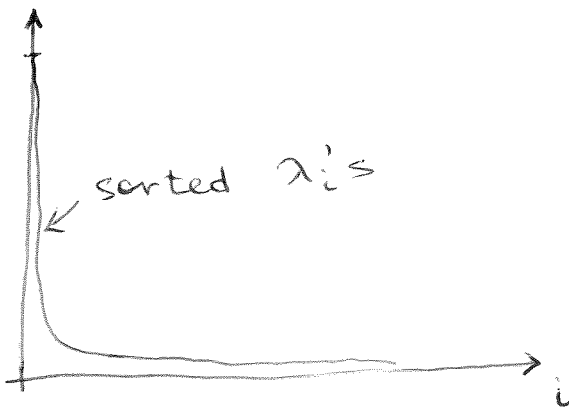
proof left as an exercise

$$S \vec{u}_i = \lambda_i \vec{u}_i \quad i=1, \dots, D$$

Then  $J = \sum_{i=M+1}^D \lambda_i$  is minimized by

choosing  $\lambda_{M+1}, \dots, \lambda_D$  to be the smallest eigenvalues  $\Rightarrow \lambda_1, \dots, \lambda_M$  are the largest eigenvalues &  $\vec{u}_1, \dots, \vec{u}_M$  define the principal subspace.

Typically, we have for high-D data:



clearly,  $J \downarrow$  as  $M \uparrow$

$$\text{Now, } \vec{x}_n \approx \sum_{i=1}^M (\vec{x}_n^T \cdot \vec{u}_i) \vec{u}_i + \sum_{i=M+1}^D (\vec{x}_n^T \cdot \vec{u}_i) \vec{u}_i \quad \ominus$$

$$\langle \vec{x} \rangle = \sum_{i=1}^D (\langle \vec{x} \rangle^T \cdot \vec{u}_i) \vec{u}_i \quad \text{gives}$$

$$\ominus \langle \vec{x} \rangle + \underbrace{\sum_{i=1}^M (\vec{x}_n^T \cdot \vec{u}_i - \langle \vec{x} \rangle^T \cdot \vec{u}_i) \vec{u}_i}_{\text{Components of an } M\text{-dim vector}}$$

This amounts to data compression

We can also use PCA to standardize the data.

Write  $S\vec{u}_i = \lambda_i \vec{u}_i$  as  $SU = UL$   
each vector is a column

$$U = \begin{pmatrix} \vec{u}_1 & \dots & \vec{u}_D \end{pmatrix} \quad L = \begin{pmatrix} \lambda_1 & & 0 \\ & \dots & \\ 0 & & \lambda_D \end{pmatrix}$$

$$U^T = U^{-1}$$

Define  $\vec{y}_n = L^{-1/2} U^T (\vec{x}_n - \langle \vec{x} \rangle)$ .

Note that  $\langle \vec{y} \rangle = \frac{1}{N} \sum_{n=1}^N \vec{y}_n = 0$ , and

covariance is then given by

$$\underbrace{\frac{1}{N} \sum_{n=1}^N \vec{y}_n \vec{y}_n^T}_{D \times D \text{ matrix}} = \frac{1}{N} \sum_{n=1}^N L^{-1/2} U^T (\vec{x}_n - \langle \vec{x} \rangle) (\vec{x}_n - \langle \vec{x} \rangle)^T * U L^{-1/2} \quad \textcircled{=}$$

$$\textcircled{=} L^{-1/2} \underbrace{U^T S U}_{U^{-1} U L = L} L^{-1/2} = \mathbb{I}.$$

This operation is known as whitening the data.

Finally,  $M=2$  projections can be used for data visualization.

# PCA for high-D data

Sometimes,  $N < D \Rightarrow$  a set of  $N$  points occupy a subspace with at most  $N-1$  dims, so using  $M > N-1$  is pointless since at least  $D - (N-1) + 1$   $\lambda_i$ 's are  $\emptyset$ . Doing "typical" PCA (i.e. finding all  $\lambda_i$ 's) scales as  $\mathcal{O}(D^3)$  and is too costly.

To address this problem, define

$$X = \begin{pmatrix} \vec{x}_1 - \langle \vec{x} \rangle \\ \vdots \\ \vec{x}_N - \langle \vec{x} \rangle \end{pmatrix}$$

Then  $\underbrace{S}_{D \times D} = \frac{1}{N} X^T X$ , and

$$\underbrace{\frac{1}{N} X^T X}_{S} \vec{u}_i = \lambda_i \vec{u}_i$$

eigenvector eq'n  $\Rightarrow \underbrace{N^{-1} X X^T}_{N \times N} (\underbrace{X \vec{u}_i}_{\vec{v}_i}) = \lambda_i (X \vec{u}_i) \quad (*)$   
 assume  $\vec{v}_i^T \cdot \vec{v}_i = 1$  w/out loss of generality

$\lambda_i, i=1 \dots N$  are the same as the "original"  $\lambda_i$ 's, and the other  $\lambda_i$ 's are all  $\emptyset$ .

Solving (\*) is an  $\mathcal{O}(N^3)$  operation.

Finally, eigenvectors:

$$\underbrace{(N^{-1} X^T X)}_S \underbrace{(X^T \vec{v}_i)}_{\text{eigenvector of } S, \text{ length } D} = \lambda_i (X^T \vec{v}_i)$$

Normalization:  $(X^T \vec{v}_i)^T (X^T \vec{v}_i) =$

$$= \vec{v}_i^T \underbrace{X X^T}_{N \lambda_i \vec{v}_i} \vec{v}_i = N \lambda_i \vec{v}_i^T \vec{v}_i$$

assuming that  $\vec{v}_i^T \vec{v}_i = 1$ , we

obtain:  $\frac{X^T \vec{v}_i}{\sqrt{N \lambda_i}}$  as normalized eigenvectors of  $S$

Thus, we construct  $N^{-1} X X^T$ , compute its eigenvalues  $\lambda_i$  & normalized eigenvectors  $\vec{v}_i$ , length  $N$

and then project back into  $D$ -space

using  $\frac{X^T \vec{v}_i}{\sqrt{N \lambda_i}}$  ← normalized eigenvectors in original data space  
length  $D$