

$$\underline{A} = \begin{bmatrix} 3 & 6 \\ 2 & 5 \\ 1 & 4 \end{bmatrix} \quad 3 \text{ rows} \quad 2 \text{ cols}$$

A is a 3 by 2 matrix

$[A]_{ij}$  : element in row  $i$  col  $j$

$L \times 1$  Let  $\underline{a}_i = \begin{bmatrix} a_{1i} \\ a_{2i} \\ \vdots \\ a_{Li} \end{bmatrix}, i = 1, 2, \dots, n$

$L \times M$   $\underline{B} = \begin{bmatrix} \underline{a}_1 & \underline{a}_2 & \dots & \underline{a}_M \\ \vdots & \vdots & & \end{bmatrix} \quad M \times L \quad \underline{C} = \begin{bmatrix} -\underline{a}_1^T \\ -\underline{a}_2^T \\ \vdots \\ -\underline{a}_M^T \end{bmatrix}$   
a as cols  
a as rows

$\underline{C} = \underline{B}^T$  : interchange rows and cols

$\underline{t}_i^T = [1 \ t_i \ t_i^2]$   $\Rightarrow$  feature associated with  $i$  response / label.

$s_i$  "label" (value goes along feature)

$$\Leftrightarrow \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{bmatrix} = \begin{bmatrix} \underline{t}_1^T \\ \vdots \\ \underline{t}_2^T \\ \vdots \\ \underline{t}_N^T \end{bmatrix} \quad w \quad \text{or} \quad \underline{s} = \begin{bmatrix} 1 & w \end{bmatrix}_{N \times 1} \quad N \times 3 \quad 3 \times 1$$

## Matrix Multiplication

$$\underline{A} : N \times M \quad \underline{B} : M \times L \quad \underline{C} : N \times L = \underline{A} \underline{B} : [C]_{m \times n} = \sum_{j=1}^m [\underline{A}]_{mj} [\underline{B}]_{jn}$$

Inner product of  $m^{th}$  row of A with  $n^{th}$  col of B

$$\underline{A} = \begin{bmatrix} 3 & 4 \\ 2 & 5 \\ 1 & 6 \end{bmatrix}, \underline{B} = \begin{bmatrix} -2 & 8 \\ 7 & -3 \end{bmatrix} \quad \underline{C} = \begin{bmatrix} 3(-2) + 4 \cdot 7 & 3 \cdot 8 + 4(-3) \\ 2(-2) + 5 \cdot 7 & 2 \cdot 8 + 5(-3) \\ 1(-2) + 6 \cdot 7 & 1 \cdot 8 + 6(-3) \end{bmatrix}$$

## Modeling multiple response

$$r: \begin{bmatrix} p_{r1} \\ p_{r2} \\ \vdots \\ p_{rL} \end{bmatrix} = \begin{bmatrix} \underline{t}_1^T \\ \vdots \\ \underline{t}_L^T \end{bmatrix} \underline{w}_r \quad b: \begin{bmatrix} p_{b1} \\ p_{b2} \\ \vdots \\ p_{bL} \end{bmatrix} = \begin{bmatrix} \underline{t}_1^T \\ \vdots \\ \underline{t}_L^T \end{bmatrix} \underline{w}_b$$

$$r+b \quad \begin{bmatrix} p_{r1} & p_{b1} \\ \vdots & \vdots \\ p_{rL} & p_{bL} \end{bmatrix} = \begin{bmatrix} \underline{w}_r & \vdots & \underline{w}_b \end{bmatrix}_{L \times 2} \quad \begin{bmatrix} \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}_{L \times 3} \quad \begin{bmatrix} \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}_{3 \times 2}$$

Multiplication rules extended to block matrices

Generalizing

$$\underline{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad \underline{B} = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

$$\underline{C} = \underline{A}\underline{B} = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}$$

All relevant submatrices must be conformable

>this: ①  $A: 5 \times 3 \quad B: 4 \times 3$  product  $AB$  is not defined.

②  $A: 5 \times 3 \quad B: 3 \times 1 \quad C = Ab : C: 5 \times 1$   
 $C$  is column vector

③  $y = Xw = \begin{bmatrix} 1x1 + 2x-1 \\ -1x1 + -1x-1 \\ 2x1 + 1x-1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$

④  $X = Tw = \begin{bmatrix} 1x1 + 2(-1) & 1x1 + 2x2 \\ -1x1 + -1(-1) & -1x1 + -1x2 \\ 2x1 + 1(-1) & 2x1 + 1x2 \end{bmatrix} = \begin{bmatrix} -1 & 5 \\ 0 & -3 \\ 1 & 4 \end{bmatrix}$

⑤  $C = [A \ b] \quad F = \begin{bmatrix} D \\ e^T \end{bmatrix} \quad A: 3 \times 2, D: 2 \times 5 \rightarrow 3 \times 5$

$$CF = \begin{bmatrix} 1 \times 2 & 2 \times 1 \\ 1 \times 1 & \end{bmatrix} \quad b: 3 \times 1 \quad e^T: 1 \times 5 \rightarrow 3 \times 5$$
$$[A \cdot D + b \cdot e^T]$$

Classification is assigning a category to data

→ extract features based on a model

if we use a line to separate classes

$$x_2 = mx_1 + b \Rightarrow x_2 - mx_1 - b = 0$$

↳ rewrite as an inner product

$$\begin{bmatrix} x_2 & x_1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \underline{x}^T \underline{w} = 0$$

↑  
features

→ classifier weights

↗ find  $\underline{w}$

↳ curved decision boundaries:

$$\begin{bmatrix} x_2 & x_1^3 & x_1^2 & x_1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_5 \end{bmatrix} = \underline{x}^T \underline{w} = 0$$

A linear classifier is based on a weighted sum of features  $\underline{x}^T \underline{w}$

Labels specify class associated with a feature binary classification (between 2)

Supervised learning: given features / labels.

$$(x_i, l_i) \text{ find } \underline{w} \text{ so } \underline{x}_i^T \underline{w} \approx l_i$$

Training a linear classifier involves solving a system of linear equations.

$$\begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \underline{w} \approx \begin{bmatrix} l_1 \\ l_2 \\ \vdots \\ l_N \end{bmatrix} \Rightarrow \begin{matrix} N \times M \\ M \times 1 \\ \uparrow \\ \underline{x}^T \underline{w} \approx \underline{l} \end{matrix}$$

$N$  training samples  
 $M$  features

Classify candidate feature  $\underline{x}$  using  $\underline{w}$

↳ if  $\underline{x}^T \underline{w} > 0 \Rightarrow$  label "+", if  $\underline{x}^T \underline{w} < 0 \Rightarrow$  label "-"

## Patterns and Model Order

$$\hat{P} = f(t) \rightarrow \text{polynomial}$$

"label" =  $[1 \ t_1 \ t_1^2]$     "feature" =  $\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$     "model" = Quadratic

building blocks-bases

$$\begin{bmatrix} \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_{20} \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \vdots & \vdots & \vdots \\ 1 & t_{20} & t_{20}^2 \end{bmatrix}}_{\mathbf{T}} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} w_0 + \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_{20} \end{bmatrix} w_1 + \begin{bmatrix} t_1^2 \\ t_2^2 \\ \vdots \\ t_{20}^2 \end{bmatrix} w_2 = \hat{P}$$

T

very simple model  $\rightarrow$  might not good in data fit

complex + data fitting model  $\rightarrow$  catching behavior not real / low generalization

- computing  $\mathbf{T}w$  involves inner product with rows of T
- Interpreting model with "bases" uses columns of T

## Modeling Matrix Data

use a "taste profile" to model each user's pref

$$\underline{t}_1 = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} \quad \sim -3\underline{t}_1 = \begin{bmatrix} -3 \\ -3 \\ 3 \end{bmatrix} \quad \sim 3\underline{t}_1 = \begin{bmatrix} 3 \\ 3 \\ -3 \end{bmatrix}$$

$\hat{R}$  can be expressed as an "outer product"

$$= \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} \underbrace{[3 \ 3 \ 3]}_{S_1^T} = \underline{t}_1 S_1^T = \begin{bmatrix} 3 & 3 & -1 & 3 \\ -3 & -3 & -1 & 3 \\ 3 & -3 & 1 & -3 \end{bmatrix}$$

$$\underline{t}_2 = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix} \quad \hat{R} = \underline{t}_1 S_1^T + \underline{t}_2 S_2^T$$

interpret matrix multiplication as sum of outer products.

$$\hat{P} = [\underline{t}_1 : \underline{t}_2] \begin{bmatrix} S_1^T \\ S_2^T \end{bmatrix} = \underline{t}_1 S_1^T + \underline{t}_2 S_2^T$$

$$\hat{R} = \underline{T} \underline{S} = [\underline{t}_1 : \underline{t}_2 : \dots : \underline{t}_L] \begin{bmatrix} S_1^T \\ S_2^T \\ \vdots \\ S_L^T \end{bmatrix} = \sum_{l=1}^L \underline{t}_l S_l^T$$

N x L    L x K

- choose L to trade model fit and generalization

- linear product for computation, outer for interpretation

$$1. C = Ab$$

$K \times 1$

$$\begin{bmatrix} a_1 & a_2 & \dots & a_m \\ a_2 \\ \vdots \\ a_k \end{bmatrix} \begin{bmatrix} b \\ \vdots \\ b_m \end{bmatrix}$$

$K \times M, \quad M \times 1 \rightarrow K \times 1$

$$2. 500 \times 10000$$

$$3. X = TW$$

$3 \times 2, 2 \times 4 \rightarrow 3 \times 4$

$$\begin{bmatrix} 2x1+2(-1) & 2x1+2x2 & 2x3+2x2 & 2x2+2x1 \\ -1(1)+1(-1) & -1(1)+1x2 & -1x3+1x2 & -1x2+1x1 \\ 2x1+1(-1) & \underline{2x1+2x1} & 2x3+1x2 & 2x2+1x1 \end{bmatrix}$$

$$4. \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \times \begin{bmatrix} 2 & -3 & -2 \end{bmatrix} \quad 3 \times 3$$

$3 \times 1$

$1 \times 3$

$$\hookrightarrow \begin{bmatrix} 1x2 & 1x(-3) & 1x(-2) \\ 2x2 & 2x(-3) & 2x(-2) \\ 3x2 & 3x(-3) & 3x(-2) \end{bmatrix}$$

$$5. X = TW \quad T = \begin{bmatrix} t_1 & t_2 \end{bmatrix}_{1 \times 2} \quad w = \begin{bmatrix} w_1^T \\ w_2^T \end{bmatrix}_{2 \times 1}$$

$$t_1 w_1^T + t_2 w_2^T$$

## 2.1. Linear dependence and Ranking in Learning

Model fitting :  $\begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_{20} \end{bmatrix} = \begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \vdots & \vdots & \vdots \\ 1 & t_{20} & t_{20}^2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$

$$\Rightarrow \underline{A} \underline{w} = \underline{d}$$

$\rightarrow$   $N \times m \times 1 \quad N \times 1$

Classifier design :  $\begin{bmatrix} l_1 \\ l_2 \\ \vdots \\ l_n \end{bmatrix} = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \underline{w}$

$N \times 1$   
↑

To solve  $\underline{A} \underline{w} = \underline{d}$ ,  $A$  can be represented as column vector  $[a_1 \dots a_m]$

$$\underline{A} \underline{w} = \underline{d} \Leftrightarrow \underline{d} = \sum_{i=1}^N a_i w_i \quad (\text{weighted sum of } a_i)$$

in graph  $\underline{d}$  is a set  $\underline{d} \neq a_1 w_1 + a_2 w_2$

$$\Leftrightarrow \underline{d} - \sum_{i=1}^m a_i w_i = 0$$

Linear independence : A set of  $M$  vectors  $v_1, v_2, \dots, v_m \in \mathbb{R}^N$  is linearly independent iff

$$\sum_{i=1}^m v_i \alpha_i = 0 \Leftrightarrow \alpha_i = 0, i=1, 2, \dots, m \quad (\text{if weights are all 0})$$

otherwise it is linearly dependent  $\leftarrow$  if there is non-zero weights that

Rank of matrix: number of linearly independent columns (or rows) (row rank = col rank)

Example :  $\underline{A} = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & -1 \\ \underline{a}_1 & \underline{a}_2 & \underline{a}_3 \end{bmatrix}$

$$a_3 = a_2 - a_1$$

$$\text{so } \underline{\alpha}_1 a_1 + \underline{\alpha}_2 a_2 + \underline{\alpha}_3 a_3 = 0$$

$$\text{for } \alpha_1 = 1, \alpha_2 = -1, \alpha_3 = 1 \neq 0$$

$a_1, a_2, a_3$  are linearly dependent

where  $a_1, a_2 \mid a_2, a_3 \mid a_1, a_3$  are independent  $\leftarrow$  can only describe  $\Rightarrow$  dimension

$\hookrightarrow$  rank of  $A$  is 2

solution :  $\underline{A} \underline{w} = \underline{d} = 0$  or

$$\sum_{i=1}^m a_i w_i + (-1)d = 0 \Rightarrow a_1, a_2, \dots, a_m, d \text{ are linearly dependent}$$

$\triangleright$  if  $d$  is a linear combination of  $a_1, a_2, \dots, a_m$  there is a solution

$$\text{rank}(A) = \text{rank}([A : d])$$

$\hookrightarrow$   $A$  with  $d$  add as an extra column  
( $d$  is dependent)

2) if  $a_1, a_2, \dots, a_m, d$  are linearly independent. There is no solution

$$\text{rank}(A) < \text{rank}([A : d])$$

Solution to  $\underline{A}\underline{w} = \underline{d}$  may not be unique

↳ Suppose  $\underline{A}\underline{w} = \underline{d}$ , does  $\underline{f} \neq 0$  exist so that  $\underline{z} = \underline{w} + \underline{f}$  also satisfy  $\underline{A}\underline{z} = \underline{d}$ ?

$$\underline{A}\underline{z} = \underline{d} \Rightarrow \underline{A}\underline{w} + \underline{A}\underline{f} = \underline{d}.$$

$$\Rightarrow (\underline{A}\underline{w} - \underline{d}) + \underline{A}\underline{f} = 0$$

we got non unique solution iff cols. in  $\underline{A}$  are linearly dependent

$$\Rightarrow \underline{A}\underline{f} = 0$$

$$\Leftrightarrow \sum_{i=1}^n a_i f_i = 0 \text{ for } \underline{f} \neq 0$$

linearly dependent

to  $\underline{A}\underline{w} = \underline{d}$ :

Solution

①  $\text{rank}\{\underline{A}\} < \text{rank}\{\underline{A} : \underline{d}\}$  no solution

②  $\text{rank}\{\underline{A}\} = \text{rank}\{\underline{A} : \underline{d}\}$

1'  $\text{rank}\{\underline{A}\} = \dim\{\underline{w}\}$  unique solution

2'  $\text{rank}\{\underline{A}\} < \dim\{\underline{w}\}$  non-unique solution

Ex:  $A: 10 \times 4 \quad \underline{w}: 4 \times 1 \quad \underline{d}: 10 \times 1$

Subspace in ML

A subspace  $S \subseteq \mathbb{R}^n$  (set of  $N$ -dim points) satisfy

1.  $0 \in S$  (contain origin)

2. if  $f, g \in S$ , then  $f+g \in S$  (close under addition)

3. if  $f \in S$ , then  $\alpha f \in S$  (close under scalar multi)

Consider  $S \subseteq \mathbb{R}^n$ ,  $\{[\underline{v}_1 \underline{v}_2 \dots \underline{v}_m] \begin{bmatrix} \underline{w}_1 \\ \underline{w}_2 \\ \vdots \\ \underline{w}_m \end{bmatrix} = \underline{V}\underline{w} \text{ for } \underline{w} \in \mathbb{R}^m\}$

$$\begin{bmatrix} \underline{w}_1 \\ \underline{w}_2 \\ \vdots \\ \underline{w}_m \end{bmatrix} \xrightarrow{\text{origin}}$$

Dimension of  $\{\underline{V}\underline{w}\}$   
is rank ( $\underline{V}$ )

1) If  $\underline{w} = 0$ ,  $\underline{V}\underline{w}$  is  $0 \in S$

2)  $f = \underline{V}\underline{w}_f \quad g = \underline{V}\underline{w}_g \quad f+g = \underline{V}(\underline{w}_f + \underline{w}_g) \in S$

3)  $f = \underline{V}\underline{w}_f \quad \alpha f = \underline{V}(\alpha \underline{w}_f) \in S$

(origin)

In general  $S = \{\underline{V}\underline{w}\}$  is a  $k = \text{rank } \underline{V} \leq \dim \text{ hyperplane in } \mathbb{R}^n$

$\underline{R} = \underline{I}S$

$\text{rank}(\underline{I}S) \leq \min\{\text{rank}(\underline{I}), \text{rank}(S)\}$

- special case:  $\underline{I} = N \times M$ ,  $\text{rank}(\underline{I}) = M$  iff  $\text{rank}(\underline{R}) = M$

$\underline{S} = M \times K$ ,  $\text{rank}(S) = M$

$\hat{\underline{R}} = \hat{\underline{R}}$ ,  $\hat{\underline{R}} = \underline{T} \underline{S} = \sum_{i=1}^m \underline{t}_i \underline{s}_i^T$  rank- $M$  approximation to  $\underline{R}$

$\hat{\underline{R}} = [\hat{\underline{t}}_1, \hat{\underline{t}}_2, \dots, \hat{\underline{t}}_k]$   $\hat{\underline{t}}_i$  lie in  $M$ -dimensional subspace

## Bases

$$S = \text{span}\{v_i\}_{i=1}^m \quad \sum_{i=1}^m v_i w_i, w_i \in \mathbb{R}$$

①  $v_i$  arbitrary - hard computing if linear dep.

②  $v_i$  linearly independent  $\sum$  Basis  $\rightarrow$  unique relationship between  $x$  and  $w_i$  in  $\Rightarrow \exists$

③  $v_i$  orthonormal - easiest computing

Orthonormal:  $v_i$  are orthogonal to each other and unit length

$$v_i^T v_i = 1, v_i^T v_j = 0 \quad \forall i \neq j$$

Examples:  $v_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, v_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$  basis for  $x-y$  plane 2D subspace in  $\mathbb{R}^3$   
 $f = v_1 w_1 + v_2 w_2$

$$\left( \begin{array}{c} \\ \vdots \\ \end{array} \right) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \begin{array}{l} v_1^T v_1 = 1 \quad v_1^T v_2 = 0 \\ \text{is orthonormal basis for } x-y \text{ plane} \end{array}$$

## Approximate solutions, Norms, and the least-squares problem

Solve exact solution for  $A \underline{w} = \underline{d}$

↳  $\sum_i a_i w_i = d$   $d$  must lie in the subspace spanned by the column of  $A$   
 hard to satisfy.

↳ can we find  $\underline{w}$  so  $A \underline{w} \approx \underline{d}$ ?  $e = A \underline{w} - \underline{d}$ ; want  $e$  small

↳ To define "small", a vector "norm" measures the size of vector

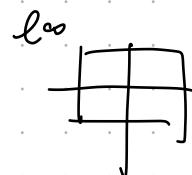
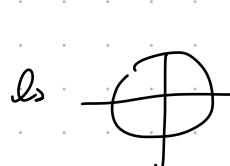
$$\|\underline{e}\|_1 = \sum_i |e_i| = a+b \rightarrow \text{Taxi-cab norm } \ell_1$$

$$\|\underline{e}\|_2 = \left( \sum_i |e_i|^2 \right)^{1/2} \rightarrow \text{Euclidean norm } \ell_2$$

$$\|\underline{e}\|_\infty = \max_i |e_i| \rightarrow \text{l-infinity norm}$$

$$\|\underline{e}\|_q = \left( \sum_i |e_i|^q \right)^{1/q} \rightarrow \text{lg norm}$$

$$\text{Unit ball: } \{x : \|x\| = 1\}$$



A vector norm  $\|\cdot\|$  maps from  $\mathbb{R}^n \rightarrow \mathbb{R}$  and satisfies:

$$\textcircled{1} \quad \|\underline{x}\| \geq 0 \text{ for all } \underline{x}$$

$$\textcircled{2} \quad \|\underline{x}\| = 0 \iff \underline{x} = 0$$

$$\textcircled{3} \quad \|b\underline{x}\| = |b| \|\underline{x}\| \text{ for all } b \in \mathbb{R}, \underline{x} \in \mathbb{R}^n$$

$$\textcircled{4} \quad \text{triangle inequality} \quad \|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\|$$

Lp norms:  $1 \leq p \leq q \leq \infty$

$$\|\underline{x}\|_q \leq \|\underline{x}\|_p$$

$$\|\underline{x}\|_\infty \leq \|\underline{x}\|$$

$\|\underline{x}\|_1 \rightarrow$  sparse solution.

$\|\underline{x}\|_\infty \rightarrow$  constant magnitude solution

$\|\underline{x}\|_2 \rightarrow$  minimize squared error

(Least-squares problem)

## The least square problem

$$\begin{array}{l} \underline{Aw} = \underline{d} \\ N \geq P \\ \text{rank } (\underline{A}) = P \end{array} \quad \min_{\underline{w}} \|\underline{Aw} - \underline{d}\|_2^2.$$

minimize

$\underline{\hat{d}} - \underline{d} \perp \text{span } (\underline{A})$

"orthogonality condition"

$$\underline{A}^T(\underline{Aw} - \underline{d}) = 0 \Rightarrow \underline{A}^T \underline{A} \underline{w} = \underline{A}^T \underline{d}$$

$$(\underline{A}^T \underline{A})^{-1} (\underline{A}^T \underline{A}) \underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d} \quad \text{matrix inverse}$$

$$\underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

## Matrix inversion

Let  $\underline{B}$  be a  $P \times P$  invertible matrix.  $\underline{B}^{-1}$  satisfies:

$$\underline{B}^{-1} \underline{B} = \underline{I} = \begin{bmatrix} 1 & & & \\ 0 & 1 & & \\ 0 & 0 & 1 & \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (\text{identity matrix } \underline{I}_v = v)$$

$$\underline{B} \underline{B}^{-1} = \underline{I}$$

$$\underline{B} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad \underline{B}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

Not all matrix have inverses. Full rank square matrix are invertible

-  $\underline{A}^T \underline{A}$  is invertible iff  $\underline{A}$  ( $N \times P$ ,  $P \leq N$ ) is rank  $P$

- positive definite  $\Rightarrow$  invertible.  $\alpha$  is positive definite ( $\alpha > 0$ ) iff  $\underline{v}^T \underline{\alpha} \underline{v} > 0 \quad \forall \underline{v} \neq 0$

$\underline{A}^T \underline{A}$  is positive definite:

Let  $y = \underline{A} \underline{v}$   $\text{rank } (\underline{A}) = P \Rightarrow y \neq 0$  for  $\underline{v} \neq 0$

$$(\underline{A} \underline{v})^T \underline{A} \underline{v} = \underline{v}^T \underline{A}^T \underline{A} \underline{v} = \underline{y}^T \underline{y} = \sum y_i^2 > 0 \quad \text{for } (\underline{v} \neq 0)$$

$\Rightarrow (\underline{A}^T \underline{A})^{-1}$  exists

## Summary

$$\min_{\underline{w}} \|\underline{Aw} - \underline{d}\|_2^2 \Rightarrow \min_{\underline{w}} \|\underline{e}\|_2^2$$

$$\Rightarrow \underline{e} \perp \text{span } \{\underline{A}\} \quad \underline{A}^T \underline{e} = \underline{0}$$

$$P \begin{bmatrix} \underline{A}^T \\ \vdots \\ \underline{A}^T \end{bmatrix} \left( \begin{bmatrix} \underline{w} \\ \vdots \\ \underline{w} \end{bmatrix} - \begin{bmatrix} \underline{d} \\ \vdots \\ \underline{d} \end{bmatrix} \right) = \underline{0}$$

$$\|\underline{e}\|_2^2 = \|\underline{d}\|_2^2 - \|\underline{\hat{d}}\|_2^2$$

$$= \underline{d}^T (\underline{I} - \underline{P}_A) \underline{d}$$

$$\underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

$$\text{projection} \quad \underline{P}_A = \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \quad \underline{\hat{d}} = \underline{P}_A \underline{d}$$

$$\underline{P}_A^* = \underline{P}_A$$

$$\underline{e} = \underline{P}_A \perp \underline{d}$$

$$\underline{P}_{A^\perp} = \underline{I} - \underline{P}_A \quad \text{projections onto space } \perp \text{ to } \text{span } \{\underline{A}\}$$

$w$  is min when  $\sum_i (y_i - \hat{y}_i)^2 = 0$

$m^T(f(w))$  is  $w = w_0$

$$= d^T d - d^T A (A^T A)^{-1} A^T d$$

$$= \|Aw - d\|_2^2$$

$\nwarrow N \text{ feature parameters} \quad \downarrow \quad \nearrow N \text{ labels}$

notes: (ignore if grading)  $[A:b]$  rank =  $[A]$

$A n = b$  if  $b \in \text{span}(\text{cols}(A))$

$\hookrightarrow A$  is linear and  $\rightarrow$  unique sol  
column space dep  $\rightarrow$  inf

$b$  is not in  $\rightarrow$  no solution

$$\begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{bmatrix}$$

feature vector      model parameter      labels

$A \underline{w} = \underline{d}$        $N \geq p$   
 $\text{rank}(A) = p$   
 $\text{full rank}$

$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2$  Let  $\hat{\underline{d}} = A\underline{w}$ ,  $\underline{d}$  lies in  $p$ -dim subspace spanned by col of  $A$

Solution  $\hat{\underline{d}} - \underline{d} \perp \text{span of col of } A \Rightarrow \underline{A}^T(\hat{\underline{d}} - \underline{d}) = 0$ .

$$\underline{A}^T(\hat{\underline{d}} - \underline{d}) = 0 \quad \text{orthogonality condition}$$

$\rightarrow$  & full rank  
 $\underline{A}^T(\underline{A}\underline{w} - \underline{d}) = 0 \Rightarrow \underline{A}^T\underline{A}\underline{w} = \underline{A}^T\underline{d}$  square matrix, have inverse.  
 $\Leftrightarrow (\underline{A}^T\underline{A})^{-1}(\underline{A}^T\underline{A})\underline{w} = (\underline{A}^T\underline{A})^{-1}\underline{A}^T\underline{d}$

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

Invertible condition:

- $A^T A$  is invertible iff  $A$  ( $N \times P$ ,  $P \leq N$ ) is rank  $P$
- Positive definite  $\Rightarrow$  invertible.  $\underline{Q}$  is positive definite ( $Q > 0$ ) iff  $\underline{v}^T \underline{Q} \underline{v} > 0 \quad \forall \underline{v} \neq 0$

$A^T A$  is positive since:

Let  $\underline{y} = A\underline{v}$  since  $A$  is full rank, and  $\underline{v} \neq 0$ ,  $A\underline{v} \neq 0$ , therefore  $\underline{y} \neq 0$

$$(A\underline{v})^T A \underline{v} = \underline{v}^T \underline{A}^T A \underline{v} = \underline{y}^T \underline{y} = \sum y_i^2 > 0 \quad (\underline{v} \neq 0)$$

$\underline{Q}$  is positive semidefinite iff  $\underline{v}^T \underline{Q} \underline{v} \geq 0 \quad \forall \underline{v} \neq 0$

$$\|\underline{z}\|_2^2 = \underline{z}^T \underline{z}$$

$$\begin{aligned} \|\underline{A}\underline{w} - \underline{d}\|_2^2 &= (\underline{A}\underline{w} - \underline{d})^T (\underline{A}\underline{w} - \underline{d}) \\ &= \underbrace{\underline{w}^T \underline{A}^T \underline{A} \underline{w} - \underline{w}^T \underline{A}^T \underline{d} - \underline{d}^T \underline{A} \underline{w} + \underline{d}^T \underline{d}}_{f(w)} \quad \leftarrow \text{we want to minimize} \end{aligned}$$

Scalar problem:  $\nabla_w f(w) = \frac{d}{dw_1} f(w) \quad \frac{d}{dw_2} f(w) \dots \frac{d}{dw_p} f(w) = 0$   $a^2 > 0 \rightarrow$  concave up

Gradients: differentiate  $f(w)$  with respect to vector  $w$   $\frac{d}{dw} f(w) = 0$

$$\nabla_w f(w) = \left[ \frac{d}{dw_1} f(w) \quad \frac{d}{dw_2} f(w) \dots \frac{d}{dw_p} f(w) \right]^T$$

Suppose  $f(w) = w^T h = h^T w = \sum_{i=1}^p h_i w_i$  with  $\frac{d}{dw_j} f(w) = h_j$

$$\nabla_w f(w) = [h_1 \ h_2 \ \dots \ h_p]^T = h$$

Suppose  $f(w) = w^T Q w$  can show  $\nabla_w f(w) = Q^T w + Qw$   
 $Q = Q^T \Rightarrow \nabla_w f(w) = 2Qw$  if  $Q$  is symmetric

and  $Q = A^T A \quad (A^T A)^T = A^T A \Leftarrow$  symmetric.

$$\nabla_w (w^T A^T A w - w^T A^T d - d^T A w + d^T d) = 0$$

$$\underline{A}^T \underline{A} \underline{w} - \underline{A}^T \underline{d} - \underline{A}^T \underline{d} = 0$$

$$2A^T A w = 2A^T d.$$

$$w = (A^T A)^{-1} A^T d$$

$$A^T A \geq 0 \Rightarrow \min f(w) \text{ when } w = w_0$$

$$\min f(w) = d^T d - d^T A (A^T A)^{-1} A^T d$$

Projection and the Pythagorean Theorem

$$\hat{d} = \underbrace{A(A^T A)^{-1} A^T d}_{P_A d} = P_A d. \quad P_A = A(A^T A)^{-1} A^T \text{ projection matrix.}$$

↳ project vector onto span plane onto A

$$P_A^2 = P_A \quad P_A^n = P_A$$

$$e = d - \hat{d} = (I - P_A)d = P_{A^\perp} d.$$

$P_{A^\perp} = I - P_A$  projects onto the space  $\perp$  to span  $\{A\}$

↳ also projection vector

$$\|d\|_2^2 = \|e\|_2^2 + \|\hat{d}\|_2^2$$

$$\|e\|_2^2 = \|d\|_2^2 - \|\hat{d}\|_2^2$$

$$= d^T d - d^T P_A d.$$

$$= d^T (I - P_A) d$$

Geometry of  $f(w)$  in  $P$ -dim space.

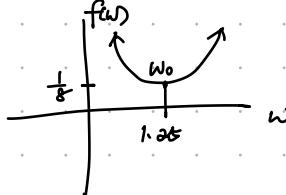
$$f(w) = (w - w_0)^T A^T A (w - w_0) + d^T P_{A^\perp} d \quad w_0 = (A^T A)^{-1} A^T d$$

$$A^T A > 0 \quad f(w) \geq f(w_0) = d^T P_{A^\perp} d$$

$$A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad d = \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} \quad w_0 = (A^T A)^{-1} A^T d = \frac{1}{2}(2, 4) = 1.25$$

$$f(w) = (w - w_0)^T A^T A (w - w_0) + d^T P_{A^\perp} d$$

$$= 2(w - w_0)^2 + \frac{1}{8}$$



Bowl shaped surface

-  $f(w) = \text{const}$  are circles

- parabolic in each coordinate ( $w$ )

if  $A^T A$  not diagonal.  $A^T A = U \Lambda^2 U^T, \Lambda^2 = \text{diag}\{\lambda_1^2, \lambda_2^2, \dots, \lambda_p^2\}$

col of  $U$  are orthonormal  $u = [u_1, u_2, \dots, u_p]$

$$u_i^T u_k = \begin{cases} 1 & i=k \\ 0 & i \neq k \end{cases}$$

$$U^T U = I$$

$$(w - w_0)^T A^T A (w - w_0) = (w - w_0)^T U \Lambda^2 U^T (w - w_0) = z^T \Lambda^2 z$$

$$w - w_0 = U z = U_1 z_1 + U_2 z_2$$

## Orthonormal basis, projections

An orthonormal basis for a set of vectors  $\mathbf{X}$  if there is another set of vectors  $\mathbf{U}$  ... such that

$$\textcircled{1} \quad u_i^T u_j = 0 \text{ for all } i \neq j$$

$$\textcircled{2} \quad u_i^T u_i = 1 \text{ for all } i$$

$$\textcircled{3} \quad \text{span}\{x_1, \dots, x_n\} = \text{span}\{u_1, \dots, u_m\}$$

subspace

$$\text{Span}\{x_1, \dots, x_n\} = \{x : x = \sum_{i=1}^n w_i x_i, w_i \in \mathbb{R}, i=1, \dots, n\}$$

$m$  is  $\dim(\text{subspace})$

properties of orthonormal basis

$$U = [u_1, \dots, u_m]$$

$$1. \quad m \leq n \quad m = n \text{ if full rank}$$

$$2. \quad m = \dim(X)$$

$$3. \quad U^T U = I_{m \times m}$$

$$4. \quad \text{if } U \text{ square } U^{-1} = U^T$$

$$5. \quad \text{if } U \text{ square } U U^T = I_{m \times m}$$

## Projections:

$$\text{proj}_{ud} = u(u^T d)$$

proj  $d$  onto  $u$       direction      amount of  $d$  in direction of  $u$

$$\text{Proj}_{ud} = \sum_{i=1}^m u_i (u_i^T d)$$

$$= [u_1, \dots, u_m] \begin{bmatrix} u_1^T d \\ \vdots \\ u_m^T d \end{bmatrix}$$

$$= U U^T d \quad \rightarrow \text{proj onto space span by } X$$

$$\text{proj}_X X = U U^T X = X$$

↓

$$U U^T = X (X^T X)^{-1} X^T \rightarrow \text{projection vector}$$

$$\text{proj}_X d = X (X^T X)^{-1} X^T d$$

• Gram-Schmidt orthogonalization

•  $\text{orth}(X)$

• SVD

Low-rank decomposition emphasize patterns.

$A(N \times M)$ ,  $T(N \times P)$ ,  $W^T(P \times M)$   $P < N$   $P \leq M$   $\text{rank}(T) = \text{rank}(W) = p$   
 $\text{rank}(TW^T) = p$

$$TW^T = \left[ \begin{array}{c|c|c|c} t_1 & t_2 & \cdots & t_p \end{array} \right] \left[ \begin{array}{c|c|c|c} \cdots & w_1^T & \cdots & \cdots \\ \cdots & w_2^T & \cdots & \cdots \\ \cdots & w_p^T & \cdots & \cdots \end{array} \right]$$
$$= \sum_{i=1}^p \underbrace{t_i w_i^T}_{\substack{\text{rank } i \\ \text{patterns}}} = \sum_{i=1}^p \underbrace{\begin{array}{c|c|c|c} \cdots & 1 & \cdots & \cdots \\ \cdots & 0 & \cdots & \cdots \\ \cdots & 0 & \cdots & \cdots \end{array}}_N \underbrace{m}_{} = \sum_{i=1}^p \boxed{n \times m}$$

↳ finding patterns

1)  $\min_{T, W} \|A - TW^T\|$  singular value decomposition

2)  $A \approx TW^T$   $T, W \geq 0$  non-negative matrix factorization

3)  $A \approx TW^T$  each col  $w_i^T$  ~~all 0 and single 1~~  $\Rightarrow$  clustering

Clustering groups similar columns.

$$\begin{bmatrix} a_1, a_2, a_3, a_4, \dots, a_m \end{bmatrix} \approx \begin{bmatrix} t_1 & t_2 & t_3 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 & \cdots & ? \\ 1 & 0 & 0 & \cdots & ? \\ 0 & 0 & 0 & \cdots & ? \end{bmatrix}$$
$$\Rightarrow a_1 \approx t_2, a_m \approx t_2, a_2 \approx t_1, \dots \quad \text{groups rows that are similar.}$$

Algorithms: K-means

so low rank models "complete" missing data.

Suppose  $a = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} w_1 + \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} w_2$

$\uparrow$  average       $\uparrow$  preference .

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

$\hat{a} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}^T \cdot \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} = \text{predicted value}$

↓ mitigate noise

Noisy data:  $\hat{A}_m = At + E$   
                 $\downarrow$  strong patterns  $\Rightarrow$  no dominant pattern.

low rank model:  $\hat{A}_m = TW^T$  can be closer to  $At$  than  $A_m$

low rank classifier / model fit

$TW^T = d$ .  $T$ : transformed features.  $x_i^T$

$$\hat{A}_m \cdot w = d \quad \left| \begin{array}{c|c|c} T & \boxed{w^T} & \\ \vdots & \downarrow & \\ w' & P \times 1 & \end{array} \right| w = \left| \begin{array}{c} d \\ \vdots \end{array} \right| \quad x_i^T = x_i^T w (w^T w)^{-1}$$

new feature  $x_i^T$ :  $x_i^T = x_i^T w (w^T w)^{-1} \Rightarrow \hat{d} = \text{sign}(x_i^T w)$

## K-means

Clustering: organizing data into groups: given  $a_i \in \mathbb{R}^N$ ,  $i=1, 2, \dots, n$ . find centroids  $\mu_j$ ,  $j=1, \dots, k$  and clusters  $S_j = \{i \mid a_i \text{ belongs to cluster } j\}$

unsupervised learning: data w/o labels, #clusters unknown

Matrix factorization:  $A \approx TW^T$   $T = [\mu_1, \mu_2, \dots, \mu_k]$

$$T = [\mu_1, \mu_2] \quad W^T = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \quad [w^T]_{lm} = \begin{cases} 1 & m \in S_l \\ 0 & m \notin S_l \end{cases}$$

K-means algorithm  $\rightarrow$  K clusters

clusters:  $S_j = \{i \mid a_i \in \text{cluster } j\}$ ,  $|S_j| = \# a_i \text{ in } S_j$

centroids:  $\mu_j = \frac{1}{|S_j|} \sum_{i \in S_j} a_i$  coherence:  $C_j = \sum_{i \in S_j} \|a_i - \mu_j\|^2$

$\hookrightarrow$  average of data vector in cluster error between data from centroid

overall coherence:  $C = \sum_{j=1}^k C_j = \sum_{j=1}^k \sum_{i \in S_j} \|a_i - \mu_j\|^2 = \|A - TW^T\|_F^2$

① Initialize: choose  $\mu_j^0$ ,  $j=1, 2, \dots, k$  randomly from  $a_{ij}$  set  $l=0$

② Assignment: put  $a_i \in S_j^l$  if  $a_i$  is closest to  $\mu_j^l$

③ update centroids  $\mu_j^{l+1} = \frac{1}{|S_j^l|} \sum_{i \in S_j^l} a_i$

④ if converged  $\rightarrow$  stop

else  $\rightarrow l = l+1$ , go to 2

options: initialization.

fixed iteration.

overall / cluster

use different norms

Challenges: convergence to local minima

unknown  $k$

SVD matrix decomposition that lead to good low rank approximations

Any  $N \times M$  matrix  $A$  can be written as  $A = U\Sigma V^T$

$N > M$

$$\begin{bmatrix} \sigma_1 & & 0 \\ 0 & \sigma_2 & \\ \vdots & & \ddots \\ 0 & & \sigma_m \end{bmatrix} \quad \downarrow \quad \begin{bmatrix} \sigma_1 & & 0 \\ 0 & \sigma_2 & \\ 0 & & \ddots \\ & & \sigma_N & 0 \end{bmatrix}$$

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_N \geq 0$  largest to smallest  
skinnier SVD

$U: N \times N$  orthonormal  
 $V: M \times M$  orthonormal  
 $\Sigma: N \times M$  diagonal  $\Sigma_{ii} \geq 0$

$$\boxed{A}_{N \times M} = \boxed{U}_{N \times N} \boxed{\Sigma}_{N \times M} \boxed{V^T}_{M \times M}$$

$\boxed{\Sigma}$  = zero out!

$$\boxed{U}_{N \times M} \boxed{\Sigma}_{M \times M} \boxed{V^T}_{M \times M}$$

$$\boxed{A}_{N \times M} = \boxed{U}_{N \times N} \boxed{\Sigma}_{N \times M} \boxed{V^T}_{M \times M} = \boxed{U}_{N \times N} \boxed{\Sigma}_{N \times N} \boxed{V^T}_{N \times M}$$

Sum of outer products form:

$$A = \begin{bmatrix} u_1 & u_2 & \dots & u_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & 0 \\ 0 & \sigma_2 & \\ & & \ddots \\ & & 0 & \sigma_m \end{bmatrix} \begin{bmatrix} -v_1^T \\ \vdots \\ -v_m^T \end{bmatrix} = \sum_{i=1}^m \sigma_i \underbrace{u_i v_i^T}_{\text{rank 1 matrix}} = \sum_{i=1}^m \sigma_i \underbrace{v_i v_i^T}_{\text{"rank 1" approximation}}$$

$$\text{Frobenius norm: } \|A\|_F^2 = \sum_{i=1}^N \sum_{j=1}^M (A)_{i,j}^2 = \|\text{vec}(A)\|_2^2$$

closest

$$\text{Eckart-Young Theorem: let } \text{rank}(A) = r \text{ and } k < r : \min_{\text{rank}(B) \leq k} \|A - B\|_F^2$$

$$= \sum_{i=k+1}^r \sigma_i^2 \text{ for } B = \sum_{i=1}^k \sigma_i u_i v_i^T \text{ where } A = U\Sigma V^T \text{ is SVD}$$

$B$  is best rank- $k$  approximation to  $A$

$\sigma_i$  provide ordered ranking of components

$$\boxed{A} \approx \sigma_1 \boxed{u_1 v_1^T} + \sigma_2 \dots + \sigma_k \dots + \sigma_{k+1} \dots$$

patterns      most important      and       $k$ -th



SVD describes matrix as an operator

$$A: N \times M \quad x: M \times 1 \quad y = Ax = U\Sigma V^T x = U \left[ \sum (\sigma_i v_i^T x) \right]$$

## Orthonormality

$$U^T U = I ; V^T V = I \rightarrow \text{full / econ}$$

$$U U^T = I_n ; V V^T = I_m \rightarrow \text{full only}$$

↪ in econ:  $U$  is  $\boxed{\quad} \quad \boxed{\quad}$   $\Rightarrow \text{rank}(U) \times \text{rank}(V)$   
cannot be  $N$  dim

$$\boxed{A_{N \times M}} = \boxed{\begin{matrix} u_1 & u_2 & \dots & u_p \\ 0 & 0 & \dots & 0 \end{matrix}} \boxed{\begin{matrix} v_1^T \\ v_2^T \\ \vdots \\ v_m^T \end{matrix}}$$

$$\boxed{A_{N \times M}} = \boxed{\begin{matrix} u_1 & u_2 & \dots & u_p \\ 0 & 0 & \dots & 0 \end{matrix}} \boxed{\begin{matrix} v_1^T \\ v_2^T \\ \vdots \\ v_m^T \end{matrix}}$$

$$\boxed{\begin{matrix} u_1 & u_2 & \dots & u_p \\ 0 & 0 & \dots & 0 \end{matrix}} \boxed{\begin{matrix} v_1^T \\ v_2^T \\ \vdots \\ v_m^T \end{matrix}}$$

## Rank

$$\text{rank}(A) = P \Leftrightarrow \sigma_1 \geq \dots \geq \sigma_p > \sigma_{p+1} = \dots = \sigma_{\min(N, M)} = 0.$$

$$A = \sum_{i=1}^P \sigma_i u_i v_i^T \quad \boxed{u} \quad \boxed{\begin{matrix} \sigma_i \\ 0 \end{matrix}} \quad \boxed{\begin{matrix} v_i^T \\ 0 \end{matrix}}$$

$$20 \times 10 \quad 10 \times 10 \quad 10 \times 10$$

$$20 \times 26$$

$$\boxed{[a_1 \ a_2 \ \dots \ a_m]} = \boxed{[u_1 \ u_2 \ \dots \ u_p]} \boxed{[c_1 \ c_2 \ \dots \ c_m]} \Rightarrow a_i = \sum_{j=1}^P u_j [c_i]_j$$

$\xrightarrow{x} \quad C = \Sigma V^T$

$a_i$  is built from  $[c_i]_j$  in basis of  $u$

left singular vector  $U$ : form an orthonormal basis from the columns of  $A$

$j^{\text{th}}$  coord  $\sim \sigma_j [c_i]_j = \sigma_j \boxed{v_i^T}_{ji}$   
 since orthonormal  
 the coordinate associated with  $j^{\text{th}}$  column of  $U$   
 is proportional to size of singular value.

right singular vector  $V$  is orthonormal basis for rows of  $A$   $D = U \Sigma$

$j^{\text{th}}$  coord  $\sim \sigma_j [d_i]_j = \sigma_j \boxed{U}_{i,j} \xrightarrow{\text{max}}$

$$N=M \quad A = U \Sigma V^T \quad U, \Sigma, V : N \times N \text{ all invertible}$$

Non invertible (singular):  $\text{rank}(A) < N$ . (not full rank)

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p > \sigma_{p+1} = \sigma_N = 0$$

Invertible:  $\text{rank}(A) = N \quad A^{-1} = V \Sigma^{-1} U^T$

$$A \cdot A^{-1} = U \Sigma V^T V \Sigma^{-1} U^T = U \Sigma \Sigma^{-1} U^T$$

$N \times N$  diagonal,  $\Sigma^{-1} = \Sigma$   
 $= U \cdot I \cdot U^T = I$ . (no econ SVD for full rank square)

$$A = \sum_{i=1}^N \sigma_i u_i v_i^T, A^{-1} = \sum_{i=1}^N \frac{1}{\sigma_i} v_i u_i^T \quad \text{SVD of } A \rightarrow \text{SVD of } A^{-1}$$

## SVD vs LSP

$$\boxed{A} = \boxed{\begin{matrix} u_1 \\ \vdots \\ u_p \end{matrix}} \quad \boxed{\begin{matrix} \Sigma \\ 0 \end{matrix}} \quad \boxed{v_1^T \dots v_m^T}$$

$$u = \boxed{\begin{matrix} \tilde{u}_1 & u_2 & \dots & u_p \end{matrix}} \quad N \times (N-P)$$

$$A = \tilde{U} \Sigma V^T$$

$$\begin{aligned} w &= (A^T A)^{-1} A^T d = (V \Sigma^2 V^T (\tilde{U} \Sigma V^T)^{-1} V \Sigma^2 \tilde{U}^T)^{-1} V \Sigma^2 \tilde{U}^T d \\ &= (V \Sigma^2 V^T)^{-1} V \Sigma^2 \tilde{U}^T d \\ &= V^T \Sigma^{-2} V^T V \Sigma^2 \tilde{U}^T d \\ &< V^T \Sigma^{-2} V^T \Sigma \tilde{U}^T d \end{aligned}$$

pseudo inverse of  $A$

$$= V^T \Sigma^{-2} \tilde{U}^T d = \sum_{i=1}^p \frac{1}{\sigma_i} V_i (\tilde{U}_i^T d)$$

$\tilde{U}$  spans  $A$

$$(E F G)^{-1} = G^{-1} F^{-1} E^{-1} \quad V^{-1} = V^T$$

$$w = V^T \Sigma^{-2} \tilde{U}^T d$$

LS error and projection

$$\hat{d} = A[(A^T A)^{-1} A^T] d = \tilde{U} \tilde{\Sigma}^{-1} \tilde{V}^T d \\ = \tilde{U} \tilde{U}^T d \quad P_A = \tilde{U} \tilde{U}^T \rightarrow \text{span}\{A\} = \text{span}\{\tilde{U}\}$$

$$e = d - \hat{d} = (I - \tilde{U} \tilde{U}^T) d = P_{A^\perp} d \quad P_{A^\perp} = I - \tilde{U} \tilde{U}^T$$

$$U^T U = I = [\tilde{U} : U_\perp] \begin{bmatrix} \tilde{U}^T \\ U_\perp^T \end{bmatrix}$$

$$I = \tilde{U} \tilde{U}^T + U_\perp U_\perp^T$$

$$P_{A^\perp} = U_\perp U_\perp^T$$

Classification  $\rightarrow$  all same norm; tell importance; can look at relative size of  $w^i$

$$A = \tilde{U} \tilde{\Sigma} V \quad \tilde{U} w^i = d \rightarrow w^i = \tilde{U}^T d \quad |w^i| \rightarrow \text{importance of } i^{\text{th}} \text{ ortho feature}$$
$$w = V \tilde{\Sigma}^{-1} w^i$$

orthobases prediction

$$\tilde{y} = \text{sign}(\tilde{x}^T w) = \text{sign}(\tilde{x}^T V \tilde{\Sigma}^{-1} w)$$

$$\tilde{x}^T = \tilde{x}^T V \tilde{\Sigma}^{-1} \quad \tilde{y} = \text{sign}(\tilde{x}^T w^i)$$

transformed feature      orthobasis classifier

$$A = U G \quad G = \tilde{\Sigma} V^T \quad x'^T = x^T c^{-1}$$

$$x'^T = x^T V^{-T} \tilde{\Sigma}^{-1}$$

$$= x^T V \tilde{\Sigma}^{-1}$$

$$x'^T = x^T V^{-T} \tilde{\Sigma}^{-T}$$

$$= x^T V^{-1} \tilde{\Sigma}^{-1}$$

$$L^2 P \quad A = U \Sigma V^T \quad w = V \Sigma^{-1} U^T d = \sum_{i=1}^p \frac{1}{\sigma_i} v_i (u_i^T d)$$

$$\Leftrightarrow \|w\|_2^2 = \sum_{i=1}^p \left(\frac{1}{\sigma_i}\right)^2 (u_i^T d)^2 \quad \text{small } \sigma_i \Rightarrow \text{large } \|w\|_2$$

$$\tilde{y} = (\bar{x} + \varepsilon)^T w = x^T w + \varepsilon^T w$$

$$|\varepsilon^T w|^2 = \|w\|_2^2 \|\varepsilon\|_2^2 \cos^2 \theta \rightarrow \text{proportional to } \|w\|_2 \quad \text{large } \|w\|_2 \rightarrow \text{amplify error}$$

$\nexists \text{ rank}(A) < p \rightarrow \sigma_p = 0 \rightarrow \text{no unique solution.}$

Regularized L2 via truncated SVD

replace  $\sum_{i=1}^p \frac{1}{\sigma_i} v_i (u_i^T d)$  with  $\sum_{i=1}^r \frac{1}{\sigma_i} v_i (u_i^T d)$  where  $r < p$

avoid small singular value.

rank-r approximation (replace  $A = \sum_{i=1}^p \sigma_i u_i v_i^T$  with  $A_r = \sum_{i=1}^r \sigma_i u_i v_i^T$ )

increase  $\min_w \|Aw - d\|_2^2$

Hidge regression

$$\min_w \|Aw - d\|_2^2 + \lambda \|w\|_2^2 \Rightarrow w = (A^T A + \lambda I)^{-1} A^T d, \quad \text{control size of norm}$$

use SVD  $A^T A = V \Sigma^2 V^T \quad \lambda I = V \lambda I V^T$

$$w = (V (\Sigma^2 + \lambda I) V^T)^{-1} V \Sigma U^T d$$

$$= V (\underbrace{\Sigma^2 + \lambda I}_{D})^{-1} \cancel{V^T} V \Sigma U^T d$$

$$= \frac{1}{\sigma_i + \lambda} v_i (u_i^T d) \quad \sigma_i \rightarrow 0 \quad \frac{\sigma_i}{\sigma_i + \lambda} \rightarrow \frac{\sigma_i}{\lambda}$$

$$w = \sum_{i=1}^p \frac{\sigma_i}{\sigma_i + \lambda} v_i (u_i^T d) \quad \text{increased value } \|Aw - d\|_2^2$$

Principle component analysis. PCA  $\Rightarrow$  maximum "variance"

Data:  $a_i, i=1, 2, \dots, n$  ( $n \times 1$ ) vectors,  $A = [a_1 \ a_2 \ \dots \ a_m]$

① assumes zero mean, center the data

$$a_i = \tilde{a}_i - \frac{1}{m} \sum_{j=1}^m \tilde{a}_j$$



PC 1: direction of accounting for maximum variance in data

$\Leftrightarrow \|f\|_2^2 = 1 \quad \max \left\{ \frac{1}{m} \sum_{i=1}^m \|a_i f\|_2^2 \right\} \text{ best line}$

unit norm

$$f_i = f^T a_i \quad \text{also } m \max_{\|f\|_2^2=1} \|a_i f\|_2^2 \Rightarrow \max_{\|f\|_2^2=1} \left\{ \underbrace{\frac{1}{m} \sum_{i=1}^m |f^T a_i|^2}_{\|f^T A\|_2^2} \right\} = \|A^T f\|_2^2$$

PC are singular vectors

$$\max_{\|f\|_2^2=1} f^T A A^T f \Rightarrow A = U \Sigma V^T \quad A A^T = U \Sigma^2 U^T$$

$$\|f\|_2^2 = 1$$

$$\max_{\|f\|_2^2=1} \frac{1}{m} f^T U \Sigma^2 U^T f \Rightarrow f = u_1$$

will maximize

Variance with w PC1  $\frac{1}{m} u_1^T A A^T u_1 = \frac{\sigma_1^2}{m}$

coord of data  $a_i = u_i^T a_i \quad a^T = [a_1 \ a_2 \ \dots \ a_m]$

$$a^T = u^T A = u^T U \Sigma V^T = \sigma_1 v_1^T$$

right singular of  $v_1$

root mean square coord:

$$\left( \frac{1}{m} \sum_{i=1}^m |\alpha_i|^2 \right)^{1/2} = \frac{1}{m^{1/2}} \|\alpha\|_2 = \frac{\sigma_1}{m^{1/2}}$$

2<sup>nd</sup> PC:  $g^T u_1 = 0$   $\|g\|_2 = 1$   $\frac{1}{m} \sum_{i=1}^m |g^T \alpha_i|^2 = g^T \alpha$   $g = v_2$  (left singular vector)

k<sup>th</sup> PC: left singular  $u_k$

$$\text{variance } \frac{\sigma_k^2}{m}$$

$$\text{root RMS } \frac{\sigma_k}{m^{1/2}}$$

rows  $\Rightarrow$  right singular  $v_i$  PC

SVD give best rank r approximation to A  $A \approx u_1 c \sigma_1 v_1^T$ ,  $r=1$

$$n = V \Sigma^{-1} U^T d$$

Data is often contaminated by noise

$$N \times M \quad A = S + C$$

M > N measured clean noise.

A: diffuse.

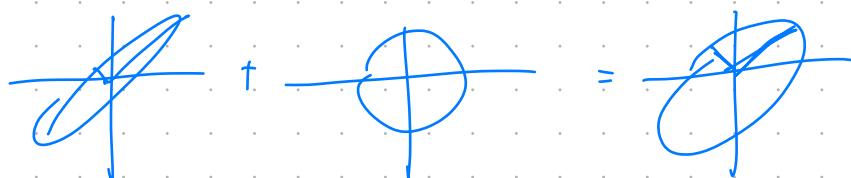
S: structured.

C: white noise - no preferred direction.

sum of squared errors:  $\|C\|_F^2$

$$\|C\|_F^2 = \sum_{i=1}^N \text{tr}\left(\frac{1}{M} \sum_{j=1}^M g_{ij}^2\right) = \text{tr}\left(\frac{1}{M} \sum_{i=1}^N \text{var}_i \approx MN \sigma_g^2\right) \quad (\text{same for all } i)$$

singular vectors are invariant (approx) to isotropic noise.



principal comp directions are unchanged.

variance along each component (singular vector vals) changes ( $M \rightarrow \infty$ )

$$U\Sigma AV^T \approx U\Sigma_S V^T + U\Sigma_C V^T$$

$$\Sigma_A = \begin{bmatrix} \sigma_{A1} \\ \sigma_{A2} \\ \vdots \\ \sigma_{AN} \end{bmatrix} \quad \Sigma_S = \begin{bmatrix} \sigma_{S1} \\ \vdots \\ \sigma_{SN} \end{bmatrix} \quad \Sigma_C \approx \begin{bmatrix} M^{\frac{1}{2}} \sigma_g \\ \vdots \\ M^{\frac{1}{2}} \sigma_g \end{bmatrix} \quad \sigma_{gi} = M^{\frac{1}{2}} \cdot \text{RMS}$$

trade bias for variance error:  $A - S \|C\|_F^2 = NM \sigma_g^2$ .

$$\text{low rank: } \hat{A}_r = \sum_{i=1}^r \sigma_{Ai} u_i v_i^T \approx \hat{S}_r + \hat{C}_r$$

$$\hat{S}_r = \sum_{i=1}^r \sigma_{Si} u_i v_i^T \quad \hat{C}_r \approx \sum_{i=r+1}^N M^{\frac{1}{2}} \sigma_g u_i v_i^T$$

$$\text{bias}^2: b^2(r) = \|S - \hat{S}_r\|_F^2$$

$$= \left\| \sum_{i=r+1}^N \sigma_{Si} u_i v_i^T \right\|_F^2$$

$$\text{variance: } V(r) = \|\hat{C}_r\|_F^2$$

$$V(r) = \left\| \sum_{i=r+1}^N M^{\frac{1}{2}} \sigma_g u_i v_i^T \right\|_F^2 = r M \sigma_g^2$$

$\sum_{i=r+1}^N \sigma_{Si}$  sum of squared "tail" singular values

↑

decrease as r increase

dimensions

variance per dimension

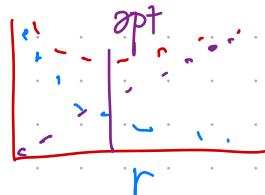
↑

increase as r increase

3. highly structured  $\sigma_1 \gg \sigma_2$ .

- $b^2(r) + V(r)$

err



- bias
- noise

weakly structured.

input bias.

Eigen decomposition applies to square matrices

Eigenvector  $e_i$ , eigenvalue  $\lambda_i$ ,  $B \in \mathbb{C}^{k \times k}$

$B e_i = \lambda_i e_i$  matrix mult  $\Leftrightarrow$  scalar mult

$$e_i \rightarrow \boxed{B} \rightarrow \lambda_i e_i$$

•  $k$  eigenvalues, possibly complex valued

• distinct  $\lambda_i \Rightarrow$  linearly independent  $e_i$

• symmetric  $B \Rightarrow$   $k$  orthonormal  $e_i$

$$E E^T = E^T E = I$$

$$\hookrightarrow B e_i = \lambda_i e_i \Rightarrow B[e_1 e_2 \dots e_k] = [e_1 e_2 \dots e_k] \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots & \\ & & & \lambda_k \end{bmatrix}$$
$$BE = E \Lambda = B = E \Lambda E^T = \sum_{i=1}^k \lambda_i e_i e_i^T$$

Symmetric PSD matrices and SVD

$$A = [a_1 \dots a_m] = U \Sigma V^T \text{ full SVD}$$

$$\begin{aligned} D B = A A^T &= \sum_{i=1}^m a_i a_i^T \\ &= U \Sigma V^T V \Sigma^T U^T = U \Sigma^2 U^T = U \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_m^2 & \\ & & 0_{00} \end{bmatrix} U^T. \end{aligned}$$

Network graph

# Unit 5

## Gradient Descent for Least Square Problems

Intuition:  $\arg \min_w \|Aw - d\|_2^2 + \gamma r(w)$

- ① Computation Cost for  $w = (A^T A)^{-1} A^T d$
- ② Closed solution might not be available
- ③ Adapt  $w$  to new feature / labels

↳ Iterative approach

$f(w) = \|Aw - d\|_2^2$        $\tau > 0$  step size

$$w^{k+1} = w^k - \tau \nabla f(w) \quad \text{as } \nabla f(w) \downarrow, \text{ step } \downarrow$$

$$f(w) = (Aw - d)^T (Aw - d)$$

$$= w^T A^T A w - 2w^T A^T d + d^T d$$

$$\nabla f(w) = 2A^T A w - 2A^T d = 2A^T(Aw - d)$$

$$w^{k+1} = w^k - \tau A^T(Aw^k - d)$$

$\tau$ : convergence behavior

$\tau$  too small: slow convergence

too big: might no convergence

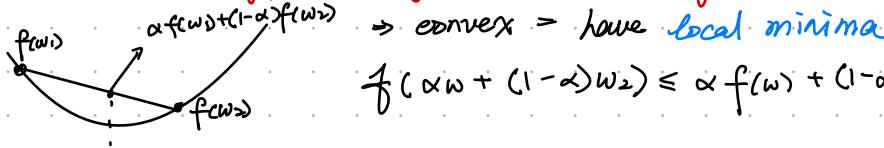
for convergence:  $f(w^{k+1}) < f(w^k) \Rightarrow$  cost decrease

$$\|Aw^{k+1} - d\|_2^2 < \|Aw^k - d\|_2^2$$

$$0 < \tau < \frac{\lambda}{\|A\|_{op}^2} = \frac{\lambda}{\sigma_{\max}^2}$$

→ guaranteed convergence to  $(A^T A)^{-1} A^T d$

Gradient descent is effective for convex cost functions



$$f(\alpha w + (1-\alpha)w_2) \leq \alpha f(w) + (1-\alpha)f(w_2) \quad \text{for } 0 < \alpha < 1$$

Multidimensional case

$$H(w) \geq 0 \quad [H(w)]_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} f(w)$$

$$\|X\|_{op} \cdot \text{largest norm } X \text{ could have} \Rightarrow \|Xw\|_2 \leq \|X\|_{op} \|w\|_2$$

Graph of contour

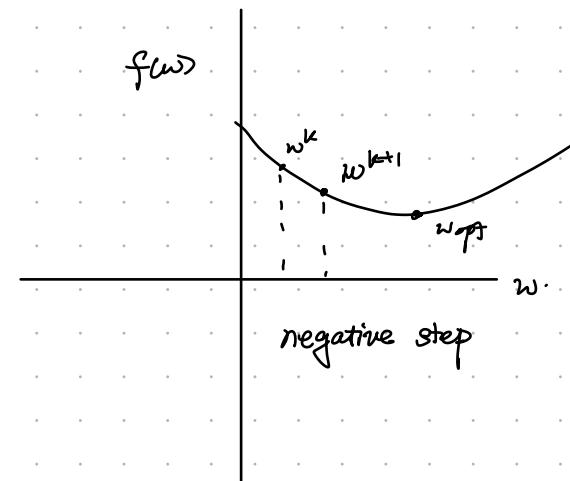
- $\sqrt{\lambda_2/\lambda_1}$  ratio ↑ the ratio between major & minor ↑
- right singular vectors define orientation of ellipse.
- if  $w$  start at major / minor axis, path is perpendicular to the  $w_L$ s
- steeper graph (ratio higher) ⇒ more steps.

Alternative way for LSP with regularizer

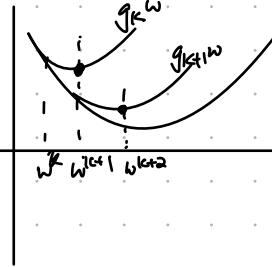
$$A^T(AA^T + \gamma I) = (A^T A + \gamma I)A^T$$

$$w = (A^T A + \gamma I)^{-1} A^T y = A^T(AA^T + \gamma I)^{-1} y$$

↳ computation cost for  $A = n \times m$



# Proximal Gradient Descent for regularized LSP



solve sequence for simpler problems

simple for separable function  $r(w) = \sum_i h_i(w_i)$

find  $g_k(w) \approx f(w) = g_k(w)$ ,  $g_k(w^k) = f(w^k)$

minimize  $g_k(w) \Rightarrow f(w)$  decrease

$$f(w) = \|d - Aw\|_2^2 + \lambda r(w)$$

$$= \| (d - Aw) + (Aw^k - Aw) \|_2^2 + \lambda r(w)$$

$$= \underbrace{\|d - Aw^k\|_2^2}_{C_k} + \underbrace{\|A(w^k - w)\|_2^2}_{\geq \|A\|_{op}^2 \|w^k - w\|_2^2} + \underbrace{2(d - Aw^k)^T A(w^k - w)}_{V_k^T} + \lambda r(w)$$

$$f(w) \leq g_k(w) = C_k + \frac{1}{2} \|w^k - w\|_2^2 + 2V_k^T(w^k - w) + \lambda r(w)$$

if  $r$  is separable function then  $g_k(w)$  is also

$$g_k(w) = C_k + \sum_i^M g_i(w_i)$$

$\hookrightarrow$  find  $\arg\min_w g_k(w) = w^{k+1}$

$$g_k(w) = C_k + \frac{1}{2} \|w^k - w\|_2^2 + 2V_k^T(w^k - w) + \lambda r(w)$$

$$\begin{aligned} t g_k(w) &= \tau C_k + (w^k - w)^T (w^k - w) + 2\tau V_k^T (w^k - w) + \lambda \tau r(w) \\ &= \tau C_k - \tau^2 V_k^T V_k + (\underbrace{\tau V_k + w^k - w}_{\zeta^{(k)}})^T (\underbrace{\tau V_k + w^k - w}_{\zeta^{(k)}}) + \lambda \tau r(w) \end{aligned}$$

$$w^{k+1} = \arg\min_w \|\zeta^k - w\|_2^2 + \lambda \tau r(w)$$

$$\zeta^k = w^k + \tau V_k$$

$$= w^k + \tau A^T (d - Aw^k)$$

$$= w^k - \tau A^T (Aw^k - d)$$

If  $r(w)$  separable  $\Rightarrow r(w) = \sum_i^M h_i(w_i)$   $w^{k+1} = \arg\min_{w_i \in \mathbb{R}} \sum_{i=1}^M (z_i^k - w_i)^2 + \lambda \tau h_i(w_i) \Rightarrow M \times \text{scalar}$

## On Ridge Regression

$$f(w) = \|d - Aw\|_2^2 + \lambda \|w\|_2^2$$

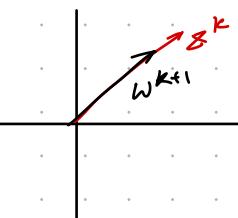
$$\text{LGD: } z^k = w^k - \tau A^T (Aw^k - d)$$

Regularization:

$$w^{k+1} = \arg\min_{w_i \in \mathbb{R}} \sum_{i=1}^M (z_i^k - w_i)^2 + \lambda \tau w_i^2$$

$$\Rightarrow w_i^{k+1} = \frac{1}{1 + \lambda \tau} z_i^k$$

$\hookrightarrow$  shrink toward origin



Graph: decrease  $\lambda \rightarrow$  less shrinkage

RR: shrinkage toward origin.

LASSO: shrinkage with offset of  $\frac{\pi}{2}$

no  $w_i w_j$  terms

$\nearrow$

$$\begin{aligned} \text{Step size or } \tau &< \frac{1}{\|A\|_{op}} \\ \frac{1}{\tau} &< \|A\|_{op} \end{aligned}$$

Alternate LS Descent  $\approx$  reg.

$$w^0 = 0, 0 < \tau < \frac{1}{\|A\|_{op}} \text{ init}$$

$$\rightarrow z^k = w^k - \tau A^T (Aw^k - d) \text{ LS GD}$$

$$w^{k+1} = \arg\min_w \|\zeta^k - w\|_2^2 + \lambda \tau r(w) \text{ reg.}$$

$\hookrightarrow$  If  $\|w^{k+1} - w^k\| < \epsilon$  stop converge

## On LASSO

$$f(w) = \|Aw - d\|_2^2 + \lambda \|w\|_1 \rightarrow \text{no closed form solution}$$

$$\text{LGD: } z^k = w^k - \tau A^T (Aw^k - d)$$

Regularization:

$$w^{k+1} = \arg\min_w \|\zeta^k - w\|_2^2 + \tau \lambda \|w\|_1$$

$$= \min_{w_i \in \mathbb{R}} \sum_{i=1}^M (z_i^k - w_i)^2 + \tau \lambda |w_i|$$

① when  $w_i > 0$ , find derivative = 0

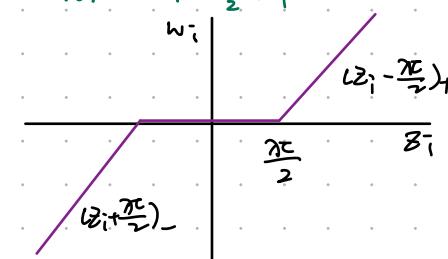
$$\min_{w_i} (z_i^k - w_i)^2 + \tau \lambda w_i$$

$$\frac{dw_i}{dz_i} \{ (z_i^k - w_i)^2 + \tau \lambda w_i \} = 0$$

$$2(z_i^k - w_i) + \tau \lambda = 0$$

$$w_i = z_i^k - \frac{\tau \lambda}{2}, w_i \geq 0$$

$$w_i = (z_i^k - \frac{\tau \lambda}{2})_+$$



$$-2(z_i^k - w_i) - \tau \lambda = 0, w_i \leq 0$$

$$w_i = z_i^k + \frac{\tau \lambda}{2}, w_i \leq 0$$

$$w_i = (z_i^k + \frac{\tau \lambda}{2})_-$$

$$w_i = \begin{cases} 0, & \frac{-\pi}{2} < z_i < \frac{\pi}{2} \\ z_i - \frac{\pi}{2}, & z_i > \frac{\pi}{2} \\ z_i + \frac{\pi}{2}, & z_i < -\frac{\pi}{2} \end{cases}$$

$$w_i = (|z_i| - \frac{\pi}{2})_+ \text{ sign}(z_i)$$

# Sparse Solutions to Least Square Problem Using LASSO

$$Aw = [a_1 \ a_2 \ \dots \ a_m] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} = \sum_{i=1}^m w_i a_i \quad \text{suppose } w_i \approx 0 \Rightarrow a_i \text{ is unimportant}$$

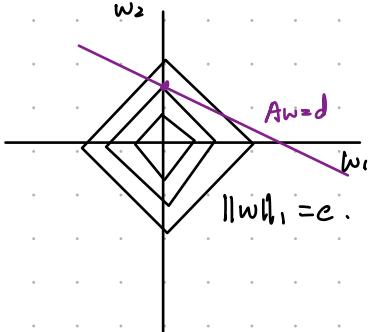
If a small number of  $w_i$  are nonzero  $\Rightarrow w$  is sparse not convex!

L1 norm:  $\|aw\|_0 + \alpha \|w\|_0 \quad \|w\|_0 = \sum_{i=1}^m \mathbb{1}_{\{w_i \neq 0\}}$  (# of non-zero elements)

**Least Absolute Selection + Shrinkage Operator is convex**

$$\min_w \|w\|_1 \text{ s.t. } \|Aw - d\|_2^2 < \epsilon$$

$$C = \sum_{i=1}^m |w_i| : |w_1| + |w_2| = C$$



$$\min_w \|w\|_1 \text{ s.t. } Aw = d$$

"corners" on  $\|w\|_1 \Rightarrow$  sparse solution

if  $\|w\|_2 \Rightarrow$  circular, non sparse

LASSO is a regularized LSP

$$\min_w \|w\|_1 \text{ s.t. } \|Aw - d\|_2^2 < \epsilon = \min_w \|Aw - d\|_2^2 + \lambda \|w\|_1 \text{ for some } \lambda, \epsilon$$

LASSO

vs

Ridge Regression

sparse  $w_L$

non sparse  $w_R$

can have small model error

can have great prediction error  $\|Aw_{opt} - Aw_R\|_2^2$

$$w_{opt} = w_L$$

can be solved in closed form

iterative solution

LASSO in model selection

$$w_L = \arg \min_w \|Aw - d\|_2^2 + \lambda \|w\|_1$$

$S_L = \{i : [w_L]_i \neq 0\}$  selected features

$$Aw_L = \sum_{i \in S_L} a_i [w_L]_i = \sum_{i \in S_L} a_i [w_e]_i$$

Debiasing

$$A_L = \{a_i : i \in S_L\}$$

$$\hat{w}_L = \arg \min_w \|A_L w - d\|_2^2$$

$$= (A_L^T A_L)^{-1} A_L^T d$$

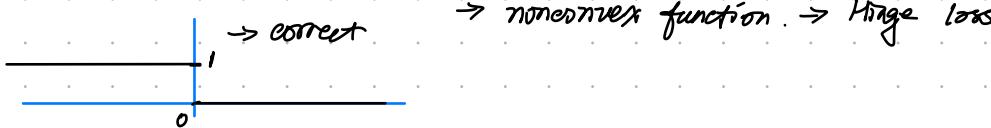
$\hookrightarrow$  avoid shrinkage due to  $\|w\|_1$

## Hinge Loss & SVM

Intuition: squared error loss is sensitive to "easy to classify" data

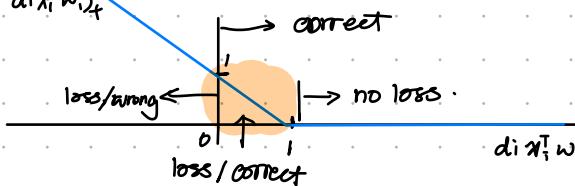
↳ result in poor classification.

Ideal loss:



Hinge loss is convex and has no loss for easy to classify data.

$$l(w; A, d) = \sum_{i=1}^N (1 - d_i x_i^T w)_+$$



Maximizing margin for separable training data.

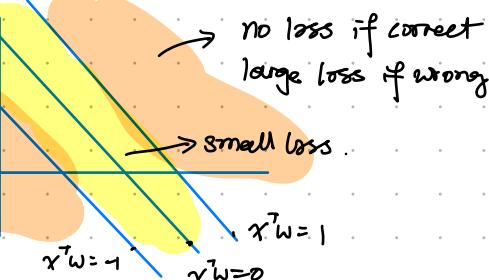
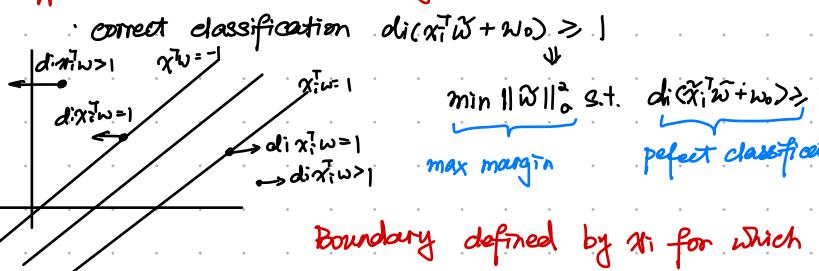
margin: distance from boundary to nearest sample

$$\text{Decision: } \hat{d} = \text{sign}(x^T w)$$

$$\hat{d} = \begin{cases} 1 & x^T w + w_0 > 0 \\ -1 & x^T w + w_0 < 0 \end{cases}$$

$$\text{boundary: } x^T w + w_0 = 0$$

Support vector maximize margin



Margin determine by  $\|w\|_2^{-1}$  ( $\frac{1}{\|w\|_2}$  distance)

unit normal to boundary plane  $v = \frac{w}{\|w\|_2}$

$$\text{Margin} = m = \frac{1}{2} \|\tilde{x}_1 - \tilde{x}_0\|_2 \quad \tilde{x}_1 = \tilde{x}_0 + 2m v$$

$$1 = \tilde{x}_1^T \tilde{w} + w_0 = \tilde{x}_0^T \tilde{w} + 2m \frac{\tilde{w}^T}{\|\tilde{w}\|_2} \tilde{w} + w_0$$

$$\text{since } \tilde{x}_0^T \tilde{w} + w_0 = -1$$

$$\geq m \frac{\tilde{w}^T}{\|\tilde{w}\|_2} \tilde{w} = 2 \rightarrow 2m \|\tilde{w}\|_2 = 2$$

$$m = \|\tilde{w}\|_2^{-1}$$

$\Rightarrow$  there is a unique solution

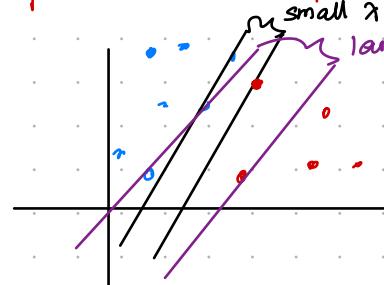
Boundary defined by  $x_i$  for which  $d_i x_i^T w = 1 \Rightarrow$  Support vectors

SVM for nonseparable data uses hinge loss

$$\min_w \sum_{i=1}^N (1 - d_i x_i^T w)_+ + \lambda \|w\|_2^2$$

loss margin norm

Only points on boundary matters.



(? future classification)

↳ when it is separable

$\hookrightarrow$  ideally, no misclassification

## Cosine Distance Descent for Support Vector Machines & Subgradients

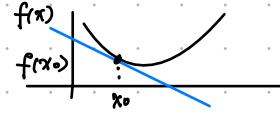
Support Vector Machine need iterative algorithms  $\min_w \sum_i^N (1 - d(x_i^T w))_+ + \gamma \|w\|_2^2$

↪ no closed form solution, but convex  $\Rightarrow$  gradient descent while hinge loss is not differentiable

Subderivatives generalize derivatives - convex, but not differentiable  $f(x)$

Normal derivative.

$$d(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0}$$



convex:  $f(x) \geq f(x_0) + d(x_0)(x - x_0)$

$$\text{Example: } f(x) = \begin{cases} \frac{1}{2}x^2 & x < 1 \\ 4x & x \geq 1 \end{cases}$$

$$d(x) = \begin{cases} x & x < 1 \\ [1, 4] & x = 1 \\ 4 & x > 1 \end{cases}$$

Subderivative (convex)

$$\text{Any } d(x_0): f(x) \geq f(x_0) + d(x_0)(x - x_0) \Rightarrow \text{blow } f(x)$$

• Gradient is generalized by Subgradients

• Gradient Descent optimization: replace gradient with subgradient

Gradient for SVMs

$$\ell(w) = \sum_i^N (1 - d(x_i^T w))_+$$

$$d(x) = \begin{cases} 1 - d(x_i^T w) & d(x_i^T w) < 1 \\ 0 & d(x_i^T w) \geq 1 \end{cases}$$

Cost:

$$f(w) = \ell(w) + \gamma \|w\|_2^2$$

$$\nabla f(w)|_{w^k} = \sum_i^N (-d(x_i) \mathbb{1}_{\{d(x_i^T w^k) < 1\}}) + 2\gamma w^k$$

Gradient Descent:

$$w^{k+1} = w^k - \tau \nabla f(w)|_{w^k}$$

Subgradient:

$$v_i(w) = \begin{cases} -d(x_i^T w) & d(x_i^T w) < 1 \\ 0 & d(x_i^T w) \geq 1 \end{cases}$$

$$= -d(x_i) \mathbb{1}_{\{d(x_i^T w) < 1\}}$$

## Stochastic Gradient Descent

• Stochastic gradient descent use part of the training data instead of all update w

Recall Squared Error

$$l(w) = \sum_{i=1}^N (d_i - x_i^T w)^2$$

$$\nabla_w l(w) = -2 \sum_{i=1}^N (d_i - x_i^T w) x_i$$

Hinge loss

$$l(w) = \sum_{i=1}^N (1 - d_i x_i^T w)_+$$

$$\nabla_w l(w) = -\sum_{i=1}^N I\{d_i x_i^T w < 1\} x_i$$

→ depend on  
all data

SGD :  $f(w) = \sum f_i(w)$  Define  $i_K$ ,  $k=1, 2, \dots$

$$w^{k+1} = w^k - \frac{\eta}{2} \nabla_w f_{i_k}(w^k)$$

SGD cycle through the training data

1) Cyclical : (incremental gradient descent)

$$i_K = k \bmod N \quad i_K = 1, 2, 3, 4, 1, 2, 3, 4, \dots$$

2) Random permutation (reshuffle every  $N$  rounds)

$$i_K = 2, 4, 1, 3, \underline{[2, 1, 4, 3]}, 4, 3, \dots$$

3) Stochastic gradient descent (uniformly at random) → most used

updated by  $-\frac{\eta}{2} \nabla_w f_{i_k}(w)$  at each iteration

$$\text{on average grad } E\{\nabla_w f_{i_k}(w)\} \approx \frac{\nabla_w f(w)}{N}$$

SGD has computational benefits

- 1) Computing  $\nabla_w f_{i_k}(w^k)$  is easier/faster than  $\nabla_w f(w^k)$
- 2) May not be able to store  $x_1, \dots, x_N$  in memory.
- 3) Noisy gradient introduces added regularization

On Ridge Regression

$\rightarrow f(w)$

$$f(w) = \sum_{i=1}^N (d_i - x_i^T w)^2 + \lambda \|w\|_2^2 = \sum_{i=1}^N (d_i - x_i^T w)^2 + \frac{\lambda}{N} \|w\|_2^2$$

$$\nabla_w f(w) = \nabla_w \left[ (d_i - x_i^T w)^2 + \frac{\lambda}{N} w^T w \right]$$

$$= -2(d_i - x_i^T w)x_i + \frac{2\lambda}{N} w$$

$$w^{k+1} = w^k - \frac{\eta}{2} \nabla_w f_{i_k}(w^k)$$

$$= w^k + 2(d_i - x_i^T w)x_i - \frac{\eta\lambda}{N} w$$

$$\text{vs } w^{k+1} = w^k + \tau A^T(Aw^k - d) - \gamma \tau w^k$$

$$A : N \times M$$

On LASSO

$\rightarrow f(w)$

$$f(w) = \sum_{i=1}^N (d_i - x_i^T w)^2 + \lambda \|w\|_1 = \sum_{i=1}^N (d_i - x_i^T w)^2 + \frac{\lambda}{N} \|w\|_1$$

Consider  $\nabla_w \sum_{i=1}^N |w_i|$

$$\frac{d}{dw_i} |w_i| = \begin{cases} \text{sign}(w_i) & w_i \neq 0 \\ [-1, 1]/0 & w_i = 0 \end{cases}$$

$$\nabla_w \|w\|_1 = \text{sign}(w)$$

$$\nabla_w f(w) = -2(d_i - x_i^T w)x_i + \frac{2}{N} \text{sign}(w)$$

$$w^{k+1} = w^k + \tau (d_{i_k} x_{i_k}^T w_{i_k} - \frac{\lambda}{2N} \text{sign}(w^k))$$

↳ for each term of  $w^k$   
take  $\text{sign}(w_i^k) \cdot 1$

## Kernel Regression

Let  $x = [x_1 \ x_2 \ \dots \ x_m]^T \in \mathbb{R}^m$

consider  $d(x) = \phi^T(x)w$ ,  $\phi(x) \in \mathbb{R}^P$   $P > m$

for example:  $x = [x_1 \ x_2]^T$

$$\phi^T(x) = [x_1^2 \ x_2^2 \ \sqrt{2}x_1x_2 \ x_1 \ 1]$$

Problem: finding  $w$  "training" data (Ridge)

$$\min_w \sum_{i=1}^N (d_i - \phi^T(x_i)w)^2 + \lambda \|w\|_2^2$$

$$\text{for } d = [d_1 \ d_2 \ \dots \ d_N]^T \quad \Phi = [\phi(x_1) \ \phi(x_2) \ \dots \ \phi(x_N)]^T \ N \times P$$

$$w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T d$$

$$d(x) = \phi^T(x)w = \phi^T(x) (\Phi^T \Phi + \lambda I)^{-1} \Phi^T d \quad \overset{P \times P}{}$$

$$= \phi^T(x) \Phi^{-1} (\Phi \Phi^T + \lambda I)^{-1} d \quad \overset{N \times N}{}$$

$$\text{since } (\Phi \Phi^T + \lambda I)^{-1} \Phi^T = \Phi^{-1} (\Phi \Phi^T + \lambda I)^{-1}$$

$$\text{Note: } [\Phi \Phi^T]_{ij} = \phi^T(x^i) \phi(x^j) \quad \left. \begin{array}{l} K(u, v) = \phi^T(u) \phi(v) \\ [\phi^T(x) \Phi^T]_j = \phi^T(x) \phi(x^j) \end{array} \right\}$$

$$\text{Let } \alpha = [\alpha_1 \ \dots \ \alpha_N]^T \quad d(x) = \sum_{i=1}^N \alpha_i \phi^T(x) \phi(x^i) = \sum_{i=1}^N \alpha_i K(x, x^i)$$

$$= (\Phi \Phi^T + \lambda I)^{-1} d$$

The Kernel to find  $d(x)$  without computing  $\phi(x)$

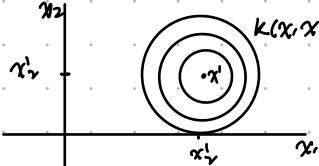
$$\text{linear: } K(u, v) = u^T v$$

$$u^T v = \|u\|_2 \|v\|_2 \cos \theta$$

$$\text{Monomials of degree } q: \phi(x) = x_1^q, x_1^{q-1} x_2, \dots, p = \frac{(q+m-1)!}{q!(m-1)!} \quad K(u, v) = (u^T v)^{\frac{q}{2}}$$

$$\text{polynomials up to degree } q: K(u, v) = (u^T v + 1)^{\frac{q}{2}}$$

$$\text{Gaussian/radial kernel: } K(u, v) = \exp \frac{-\|u-v\|^2}{2\sigma^2}$$



- no explicit  $\phi(x)$
- all polynomial orders
- smoothness controlled by  $\sigma$

## Kernel regression considerations

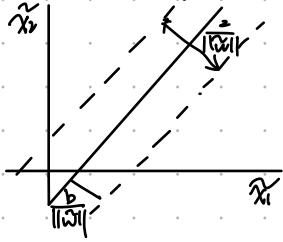
$$d(x) = \phi^T(x)w \quad \text{vs} \quad d(x) = \sum_{i=1}^N \alpha_i K(x, x^i)$$

① Store and compute  $\alpha (N \times 1)$  vs  $w (P \times 1)$

② Binary classification  $\text{sign}(d(x))$

③ Avoid "overfitting" with high dim feature space  $\Rightarrow$  use cross validation.

## Kernel Based SVM



$w = [\tilde{w}^\top \ b]$ , depends only on support vectors

$$x^\top w = 0$$

$$d(x) = \phi^\top(\alpha) w = \sum_{j=1}^N \alpha_j k(x, x^j)$$

$$\text{Let } \phi(x) = x \Rightarrow k(x, x^j) = x^\top x^j$$

$$d(x) = \sum_{j=1}^N \alpha_j x^\top x^j = x^\top \sum_{j=1}^N \alpha_j x^j$$

$$w = \sum_{j=1}^N \alpha_j x^j \Rightarrow \text{all } \alpha_j = 0 \text{ except support vectors}$$

High-dimensional feature space

$$\hat{d}(x) = \text{sign}(\phi^\top(\alpha) w)$$

Hinge loss with ridge regression

$$\min_w \sum_{i=1}^N (1 - d_i \phi^\top(\alpha) w)_+ + \lambda \|w\|_2^2 \Rightarrow x = \sum_{j=1}^N \phi(x^j) \alpha_j$$

Kernel trick replace  $\phi^\top(\alpha) \phi(x^j)$  with  $K(x^i, x^j)$

$$\begin{aligned} & \min_{\alpha} \sum_{i=1}^N (1 - d_i \phi^\top(\alpha)) \underbrace{\sum_{j=1}^N \alpha_j \phi(x^j)}_w)_+ + \lambda \sum_{i=1}^N \alpha_i \phi^\top(\alpha) \sum_{j=1}^N \alpha_j \phi(x^j) \\ &= \min_{\alpha} \sum_{i=1}^N (1 - d_i \sum_{j=1}^N \alpha_j K(x^i, x^j))_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(x^i, x^j) \\ & \quad \text{sum: } \min_{\alpha} \sum_{i=1}^N (1 - d_i \sum_{j=1}^N \alpha_j K(x^i, x^j))_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(x^i, x^j) \end{aligned}$$

Support Vector Machine has sparse  $\alpha$

$$\text{decision boundary: } d(x) = 0 = \phi^\top(\alpha) w = \sum_{j=1}^N \alpha_j K(x, x^j)$$

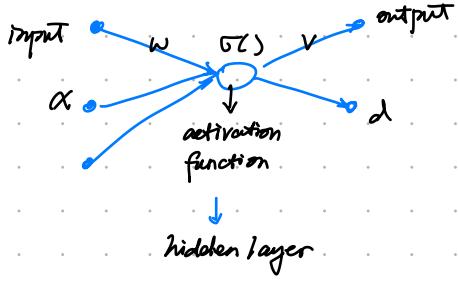
- Boundary only depends on support vectors
- $K(u, v)$  measures the similarity of  $u, v$
- solve for  $\alpha$  use gradient descent.

Graphing:

- larger  $\lambda \Rightarrow$  smaller kernel norm (larger penalty)  $\Rightarrow$  smooth, underestimation
  - larger  $\sigma \Rightarrow$  wider kernel, more smooth,  
really small  $\sigma \Rightarrow$  reply to noise,
- $\hookrightarrow$  if in extreme case width  $\Rightarrow 0$ , prediction only depends on closest sample.
- when kernel width grow  $\Rightarrow$  account for more training samples for prediction

a Kernel have maximum value when  $x_i = x^j$  for  $K(x, x^j)$   
 $x^j$  is usually center of kernel

## Neural Network



all  $\circ$  : neuron

$$\hat{d} = \sigma(x^T w)$$

Common  $\sigma(z)$

ReLU

$$\sigma(z) = \max(0, z)$$

logistic

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

Sign

$$\sigma(z) = \text{sign}(z)$$

multiple outputs  $\Rightarrow$  solve multiple problems

$$h_k = \sigma\left(\sum_j w_{kj} x_j\right)$$

$$\hat{d}_m = \sigma\left(\sum_n v_{mn} h_n\right)$$

pending Problems

• Network structure : # of layers, # of hidden nodes

• Choose the weights  
 $\hookrightarrow$  SGD and backpropagation (Nonconvex?)

Backpropagation uses SGD to train weights.

- $N$  training samples  $(x_1^i, x_2^i \dots x_m^i, d^i; i=1, 2 \dots N)$

1) Initialize  $w_{k,l}, v_{k,n}$

2) for  $t=1, 2, 3 \dots$

pick  $i_t \in \{1, 2, 3 \dots, N\}$  randomly, Compute  $h_j^{it}$ ,  $\hat{d}_k^{it}$ ,

$$\text{for } \frac{\partial f^i}{\partial v_{k,l}} : \left(\frac{\partial f^i}{\partial \hat{d}_k}\right) \frac{\partial \hat{d}_k}{\partial v_{k,l}}$$

$$\frac{\partial f^i}{\partial v_{k,l}} = (\hat{d}_k^{it} - d_k^i) \cdot \sigma'(\sum_n h_n^i v_{kn}) \cdot \frac{\partial}{\partial v_{k,l}} (\sum_n h_n^i v_{kn})$$

$$= \underbrace{(\hat{d}_k^{it} - d_k^i)}_{\delta_k^{it}} \underbrace{\hat{d}_k^{it} (1 - \hat{d}_k^{it})}_{\text{layer input}} h_e^i$$

$$= \delta_k^{it} h_e^i$$

$$v_{k,l}^{t+1} = v_{k,l}^t - \alpha \delta_k^{it} h_e^i$$

$$v_{k,l}^{t+1} = v_{k,l}^t - \alpha \frac{\partial f^i}{\partial v_{k,l}}$$

$$w_{m,j}^{t+1} = w_{m,j}^t - \alpha \frac{\partial f^i}{\partial w_{m,j}}$$

Chain rule

$$\text{for } \sigma(z) = \frac{1}{1+e^{-z}}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\sigma(\sum_n h_n^i v_{kn})(1 - \sigma(\sum_n h_n^i v_{kn}))$$

$$= \hat{d}_k^{it} (1 - \hat{d}_k^{it})$$

$$\text{for } \frac{\partial f^i}{\partial w_{m,j}} : \sum_q \frac{\partial f^i}{\partial \hat{d}_q} \cdot \frac{\partial \hat{d}_q}{\partial h_m^i} \cdot \frac{\partial h_m^i}{\partial w_{m,j}}$$

$$\frac{\partial f^i}{\partial w_{m,j}} = (\hat{d}_q^{it} - d_q^i) \hat{d}_q^{it} (1 - \hat{d}_q^{it}) V_{q,m}^i h_m^i (1 - h_m^i) x_j^i = \gamma_m^i x_j^i$$

$$w_{m,j}^{t+1} = w_{m,j}^t - \alpha \gamma_m^i x_j^i$$

## Summary

Initialize :  $w_{m,j}^0 > v_{k,l}^0$

Iterate for  $t=0, 1, 2 \dots$

SGD choose  $i_t \in \{1, 2, \dots, N\}$  at random

Forward Net compute  $h_m^i$ ,  $\hat{d}_q^{it}$  from  $x_e^i$ ,  $w_{m,j}^t$ ,  $v_{k,l}^t$

Backward update  $v_{k,l}^{t+1} = v_{k,l}^t - \alpha \delta_k^{it} h_e^i$

$$\gamma_m^i = \sum_q \delta_q^{it} V_{q,m}^i h_m^i (1 - h_m^i)$$

$$w_{m,j}^{t+1} = w_{m,j}^t - \alpha \gamma_m^i x_j^i$$