



Networked Life

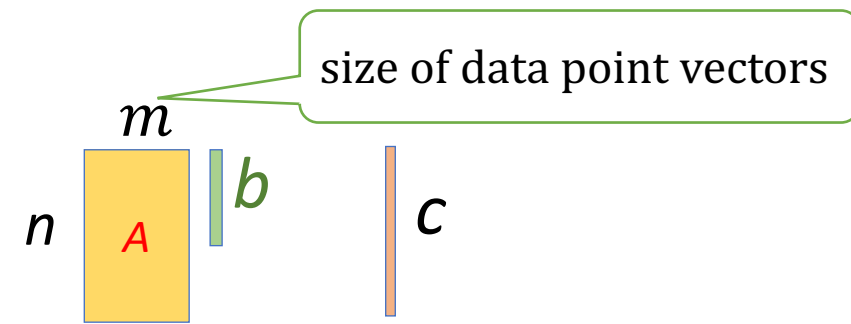
## Q4 Tutorial

- **Matrix calculus**
- **Applications**

# Optimization

- We want to solve  $\min_b \|Ab - c\|_2$
- This is a convex problem (quadratic in  $b$ )
- solve for  $b$  the **consistent** set of equations  $(A^T A)b = A^T c$
- If  $A$  has independent columns, then  $A^T A$  is invertible and there is a unique solution  $b^* = (A^T A)^{-1} A^T c$

# of data points



We provide two proofs in the next slides

- a) by linear algebra
- b) by matrix calculus

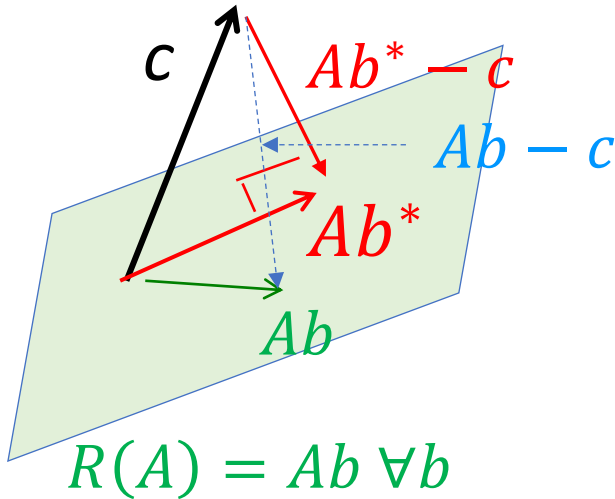
Diagram illustrating the multiplication of  $A^T$  and  $A$  to form  $A^T A$ :

$$\begin{bmatrix} A^T \end{bmatrix} \begin{bmatrix} A \end{bmatrix} = \begin{bmatrix} A^T A \end{bmatrix}$$

# Proof A

(\*), (\*\*) based on [https://ocw.mit.edu/courses/mathematics/18-06sc-linear-algebra-fall-2011/least-squares-determinants-and-eigenvalues/orthogonal-vectors-and-subspaces/MIT18\\_06SCF11\\_Ses2.1sum.pdf](https://ocw.mit.edu/courses/mathematics/18-06sc-linear-algebra-fall-2011/least-squares-determinants-and-eigenvalues/orthogonal-vectors-and-subspaces/MIT18_06SCF11_Ses2.1sum.pdf)

- $\min_b \|Ab - c\|_2$  is equivalent to  $(A^T A)b = A^T c$



Minimize the distance of  $c$  and  $Ab$  over all possible  $b$ s:  
project  $c$  on the linear subspace range of  $A = R(A) \Leftrightarrow$   
find  $b^*$  s. t.  $Ab^* - c \perp R(A)$   
 $\Leftrightarrow Ab^* - c \in N(A^T)$  (\*)  
 $\Leftrightarrow A^T(Ab^* - c) = 0$

- If  $A$  has independent columns, then  $A^T A$  is invertible and there is a unique solution  $b^* = (A^T A)^{-1} A^T c$  (\*\*)

$$\begin{array}{|c|} \hline A^T \\ \hline \end{array} \begin{array}{|c|} \hline A \\ \hline \end{array} = \begin{array}{|c|} \hline A^T A \\ \hline \end{array}$$

# Proof B: matrix calculus

# The key matrix calculus properties

scalars = italics

$$a = \mathbf{y}^T \mathbf{x} = \mathbf{x}^T \mathbf{y} \Leftrightarrow \frac{da}{d\mathbf{x}} = \left( \frac{da}{dx_1}, \dots, \frac{da}{dx_n} \right) = \mathbf{y}^T \quad (*)$$

$$\mathbf{y} = \mathbf{B}\mathbf{x} \Leftrightarrow \frac{d\mathbf{y}}{d\mathbf{x}} = \frac{d}{d\mathbf{x}}(y_1(\mathbf{x}), \dots, y_n(\mathbf{x}))^T = \mathbf{B} \quad (**)$$

$$a = \mathbf{y}(\mathbf{x})^T \mathbf{z}(\mathbf{x}) \Leftrightarrow \frac{da}{d\mathbf{x}} = \mathbf{y}(\mathbf{x})^T \frac{d\mathbf{z}(\mathbf{x})}{d\mathbf{x}} + \mathbf{z}(\mathbf{x})^T \frac{d\mathbf{y}(\mathbf{x})}{d\mathbf{x}} \quad (***)$$

$$a = \mathbf{x}^T \mathbf{B}\mathbf{x} \Leftrightarrow \frac{da}{d\mathbf{x}} \stackrel{(***)}{=} \mathbf{x}^T \frac{d(\mathbf{B}\mathbf{x})}{d\mathbf{x}} + (\mathbf{B}\mathbf{x})^T \frac{d\mathbf{x}}{d\mathbf{x}} \stackrel{(**)}{=} \mathbf{x}^T (\mathbf{B} + \mathbf{B}^T) \quad (1)$$

For matrix calculus please see section 5 in

<http://www.atmos.washington.edu/~dennis/MatrixCalculus.pdf>

# Minimizing SE (calculation to practice matrix calculus)

$$\min_b ||Ab - y||^2 = \min_b (Ab - y)^T (Ab - y)$$

Hence need to find  $b$  s.t.  $\frac{d}{db} ||Ab - y||^2 = 0$

$$||Ab - y||^2 = (Ab - y)^T (Ab - y) = b^T A^T Ab - y^T Ab - b^T A^T y + y^T y$$

$$\frac{d}{db} ||Ab - y||^2 = b^T \frac{d}{db} (A^T Ab) + (A^T Ab)^T \frac{d}{db} b - y^T A \frac{d}{db} b - (A^T y)^T \frac{d}{db} b + 0$$

$$= b^T A^T A + b^T A^T A - y^T A - y^T A$$

$$= 2(b^T A^T - y^T)A = 2(Ab - y)^T A$$

$$\text{Hence } \frac{d}{db} ||Ab - y||^2 = 0 \Leftrightarrow (Ab - y)^T A = 0$$

$$\Leftrightarrow A^T (Ab - y) = 0 \Leftrightarrow A^T Ab = A^T y$$

# Minimizing SE (simpler derivation)

$$\min_b ||Ab - y||^2 = \min_b (Ab - y)^T (Ab - y)$$

Hence need to find  $b$  s.t.  $\frac{d}{db} (Ab - y)^T (Ab - y) = 0$

$$\begin{aligned} \frac{d}{db} (Ab - y)^T (Ab - y) &= (Ab - y)^T \frac{d}{db} (Ab - y) + (Ab - y)^T \frac{d}{db} (Ab - y) \\ &= 2(Ab - y)^T A \end{aligned}$$

$$\begin{aligned} \frac{d}{db} ||Ab - y||^2 = 0 &\Leftrightarrow 2(Ab - y)^T A = 0 \\ &\Leftrightarrow A^T (Ab - y) = 0 \\ &\Leftrightarrow A^T Ab = A^T y \end{aligned}$$

# Latent factor example (next 2 slides)

## 1. Small data set case: over fitting

- Number of variables (6) > number of ratings (5)
- We can match the training set exactly!

## 2. Larger data set case

- Number of variables (8) = number of ratings (8)
- Cannot match exactly the ratings since it would require some of the variables to be  $< 0$ . Hence, we do not over fit (we see better results for our forecasted ratings for the test set)