

CEE 616: Probabilistic Machine Learning

Lecture 1e: Foundations—Linear Algebra

Jimi Oke

UMass **Amherst**
College of Engineering

Feb 16, 2025

Outline

- 1 Introduction
- 2 Vectors
- 3 Matrices
- 4 Special matrices
- 5 EVD
- 6 Linear systems
- 7 Outlook

Scalars, vectors and matrices

- **Scalar:** a single number, e.g. $a \in \mathbb{R}$
- **Vector:** ordered array of numbers, e.g. $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

- **Matrix:** two-dimensional array of numbers, e.g. $\mathbf{X} \in \mathbb{R}^{n \times p}$

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

Each element $x_{ij} = [\mathbf{X}]_{ij} \in \mathbb{R}, \forall i \in \{1 : n\}, j \in \{1 : p\}$

Tensors

- **Tensor:** generalization of a matrix to arbitrary number of indices/dimensions, e.g. 3 (x_{ijk})
- Number of dimensions is called **order/rank**
- A common application is the representation of an RGB image, e.g. a square 256-pixel image can be denoted by $\mathbf{A} \in \mathbb{R}^{256 \times 256 \times 3}$

Vector and matrix multiplication

Dot/inner product of two vectors

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^\top \mathbf{y} = \begin{bmatrix} x_1 & \cdots & x_N \end{bmatrix} \times \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \sum_{n=1}^N x_n y_n \quad (1)$$

Matrix product

$$\mathbf{A}_{M \times N} \mathbf{B}_{N \times D} = \mathbf{C}_{M \times D} \quad (2)$$

Each element $[\mathbf{C}]_{md}$ is obtained as the dot product between the m th row of \mathbf{A} and the d -th column of \mathbf{B} .

Matrix multiplication properties

- Distributivity

$$\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC} \quad (3)$$

- Associativity

$$\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C} \quad (4)$$

- Conjugate transposability

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top \quad (5)$$

Non-commutativity

Unlike the inner product of two vectors, where $\mathbf{x}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{x}$, matrix multiplication is not commutative:

$$\mathbf{AB} \neq \mathbf{BA} \quad (6)$$

Linear independence

A set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is said to be linearly independent (LI) if no vector in the set can be expressed as a linear combination of the others.

- If a vector \mathbf{x}_j in the set can be written as:

$$\mathbf{x}_j = \sum_{i \neq j}^n \alpha_i \mathbf{x}_i \quad (7)$$

Then \mathbf{x}_j is dependent on the others, so the set is not LI

- Another way to view this:

$$\nexists \boldsymbol{\alpha} \in \mathbb{R}^n \quad \text{s.t.} \quad \alpha_1 \mathbf{x}_1 + \dots + \alpha_n \mathbf{x}_n = \mathbf{0} \quad (8)$$

Span

The **span** of a set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is the set of all vectors that can be expressed as a linear combination of the set:

$$\text{span}(\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}) := \left\{ \mathbf{v} : \mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{x}_i, \quad \alpha_i \in \mathbb{R} \right\} \quad (9)$$

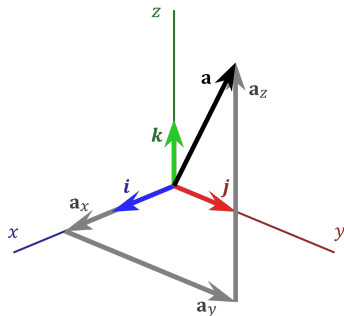
Result

If a set of vectors is LI, then any vector $\mathbf{v} \in \mathbb{R}^n$ can be written as a linear combination of \mathbf{x}_1 through \mathbf{x}_n

Basis

A set of LI vectors that span an entire vector space \mathcal{V} is called a **basis** \mathcal{B} .

- Every element of $\mathbf{v} \in \mathcal{V}$ can be expressed as a linear combination of corresponding elements in \mathcal{B} .



Standard basis vectors in \mathbb{R}^3 : $\mathbf{i}, \mathbf{j}, \mathbf{k}$ or $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$.

Source: https://en.wikipedia.org/wiki/Standard_basis

Vector norms

Norm of a vector (p-norm)

$$||\mathbf{x}||_D = \left(\sum_n |x_n|^D \right)^{\frac{1}{D}} \quad (10)$$

- Euclidean, ℓ^2 , or 2-norm:

$$||\mathbf{x}||_2 = \sqrt{\sum_{n=1}^N x_n^2} \quad (11)$$

$$||\mathbf{x}||_2^2 = \mathbf{x}^\top \mathbf{x}$$

- ℓ^1 or 1-norm:

$$||\mathbf{x}||_1 = \sum_{n=1}^N |x_n| \quad (12)$$

- Max-norm:

$$||\mathbf{x}||_\infty = \max_n |x_n| \quad (13)$$

- Unit vector: one whose Euclidean norm is 1: $||\mathbf{x}||_2 = 1$

Special vectors

- The vector of ones: **1**.
- Vector of zeros: **0**.
- Unit (or one-hot) vector: all zeros except for entry i with a value of 1.

$$\mathbf{e}_4 \in \mathbb{R}^5 = (0, 0, 0, 1, 0) \quad (14)$$

Transpose and trace

- The transpose \mathbf{A}^\top of a matrix \mathbf{A} is obtained by converting row elements to column elements and vice versa:

$$[\mathbf{A}^\top]_{ij} = [\mathbf{A}]_{ji} \quad (15)$$

If \mathbf{A} is an $n \times p$ matrix, then \mathbf{A}^\top is $p \times n$.

- A matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ is **square** if $n = p$.
- If a square matrix \mathbf{A} is equal to its transpose ($\mathbf{A} = \mathbf{A}^\top$), then \mathbf{A} is **symmetric**.
- The trace $\text{tr}(\mathbf{A})$ of a square matrix \mathbf{A} is the sum of its diagonal elements:

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii} \quad (16)$$

Matrix norms

Given a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ defining a function $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$, the **induced p-norm** of \mathbf{A} is given by:

$$\|\mathbf{A}\|_D = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|_D \quad (17)$$

For $D = 2$:

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} = \max_n \sigma_n \quad (18)$$

where λ_{\max} is the greatest eigenvalue and σ_n is n -th singular value.

Determinants

Geometrically, the determinant, $\det(\mathbf{A})$ or $|\mathbf{A}|$, of a square matrix is the [directional] scaling factor of a unit area/volume transformed by the matrix.

- If $\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$, then:

$$|\mathbf{A}| = ad - bc \quad (19)$$

- The determinant of a singular (noninvertible) matrix is 0.
- The determinant of matrix is equal to the product of its eigenvalues:

$$|\mathbf{A}| = \prod_{i=1}^n \lambda_i \quad (20)$$

Range and nullspace

The range (column space) of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the set of all the vectors that can be expressed as a linear combination of the column vectors of \mathbf{A} .

$$\text{range}(\mathbf{A}) := \{\mathbf{v} \in \mathbb{R}^m : \mathbf{v} = \mathbf{A}\mathbf{x}, \mathbf{x} \in \mathbb{R}^n\} \quad (21)$$

The nullspace of a matrix is the set of all vectors \mathbf{v} such that $\mathbf{A}\mathbf{v} = \mathbf{0}$:

$$\text{nullspace}(\mathbf{A}) := \{\mathbf{v} \in \mathbb{R}^n : \mathbf{A}\mathbf{v} = \mathbf{0}\} \quad (22)$$

Rank

The rank of a matrix is the greatest number of its LI column vectors (or row vectors)

- This is equivalent to the dimension of the vector space spanned by its columns (or by its rows)
- Given $\mathbf{A} \in \mathbb{R}^{m \times n}$:

$$\text{rank}(\mathbf{A}) \leq \min(m, n) \quad (23)$$

- A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is full rank if $\text{rank}(\mathbf{A}) = \min(m, n)$.
- Any matrix that is not full rank is **rank deficient**
- A square matrix is invertible if and only if it is full rank
- $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^\top) = \text{rank}(\mathbf{A}^\top \mathbf{A}) = \text{rank}(\mathbf{A} \mathbf{A}^\top)$

Condition number

The condition number κ of a matrix measures how numerically stable it is under computation:

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| \quad (24)$$

where $\|\mathbf{A}\|$ is typically taken as the ℓ_2 norm of \mathbf{A} .

- \mathbf{A} is well-conditioned when $\kappa(\mathbf{A})$ is close to 1
- \mathbf{A} is ill-conditioned when $\kappa(\mathbf{A})$ is large (nearly singular/noninvertible)

Special matrices

- **Diagonal matrix:** square matrix; all elements 0 except on the main diagonal

$$\mathbf{A} = \text{diag}(\mathbf{a}) = \begin{bmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_n \end{bmatrix}$$

- **Identity matrix:** diagonal matrix whose non-zero elements are 1

$$\mathbf{I}_n = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}_{n \times n}$$

- **Block diagonal matrix:** concatenates matrices onto the main diagonal of a single one

$$\mathbf{Z} = \text{blkdiag}(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \begin{bmatrix} \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{C} \end{bmatrix}$$

Triangular matrices

- Upper triangular matrix: all non-zero entries are either on or above the diagonal:

$$\begin{pmatrix} 1 & 4 & 7 & -2 \\ 0 & -3 & 1 & 10 \\ 0 & 0 & 7 & 2 \\ 0 & 0 & 0 & -1 \end{pmatrix}$$

- Lower triangular matrix: all non-zero entries are either on or below the diagonal:

$$\begin{pmatrix} 3 & 0 & 0 \\ -1 & -2 & 0 \\ 9 & 4 & 1 \end{pmatrix}$$

- The diagonal elements A_{ii} of a triangular matrix \mathbf{A} are its eigenvalues.

Positive definite matrices

Given a symmetric matrix $\mathbf{A} \in \mathbb{S}^n$ and a non-zero vector $\mathbf{x} \in \mathbb{R}^n$, \mathbf{A} is:

- positive definite if $\forall \mathbf{x}: \mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$
- positive semidefinite (psd) if $\forall \mathbf{x}: \mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$
- negative definite if $\forall \mathbf{x}: \mathbf{x}^\top \mathbf{A} \mathbf{x} < 0$
- negative semidefinite if $\forall \mathbf{x}: \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq 0$
- indefinite if neither psd nor nsd
- **Gram matrix \mathbf{G}** is always psd:

$$\mathbf{G} = \mathbf{X}^\top \mathbf{X} \quad (25)$$

given any matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$

Orthogonal matrix

- Two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ are orthogonal if $\mathbf{x}^\top \mathbf{y} = 0$.
- A normalized vector is one whose 2-norm is 1.
- Thus, a set of vectors that is pairwise orthogonal and normalized is **orthonormal**
- A square matrix $\mathbf{U} \in \mathbb{R}^{n \times n}$ is **orthogonal** if all its columns are orthonormal
- The inverse of \mathbf{U} is its transpose:

$$\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I} \quad (26)$$

which implies: $\mathbf{U}^{-1} = \mathbf{U}^\top$

Nonsingular matrices

A matrix is nonsingular (invertible) only if it is square and if its columns are linearly independent.

The inverse of a matrix \mathbf{X} is given by \mathbf{X}^{-1} and satisfies the property:

$$\mathbf{X}^{-1}\mathbf{X} = \mathbf{X}\mathbf{X}^{-1} = \mathbf{I}_n \quad (27)$$

- \mathbf{X}^{-1} exists $\iff |\mathbf{X}| \neq 0$
- $(\mathbf{X}^{-1})^{-1} = \mathbf{X}$
- $(\mathbf{W}\mathbf{X})^{-1} = \mathbf{X}^{-1}\mathbf{W}^{-1}$
- For 2D case:

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, \quad \mathbf{A}^{-1} = \frac{1}{|\mathbf{A}|} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \quad (28)$$

- For block diagonal matrix:

$$\begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{B}^{-1} \end{pmatrix} \quad (29)$$

Other important results

Read **PMLI** 7.3 for details

- Schur complement: Given a partitioned matrix $\mathbf{M} = \begin{pmatrix} \mathbf{E} & \mathbf{F} \\ \mathbf{G} & \mathbf{H} \end{pmatrix}$ then:

$$\mathbf{M}/\mathbf{H} = \mathbf{E} - \mathbf{F}\mathbf{H}^{-1}\mathbf{G} \quad (\text{Schur complement of } \mathbf{M} \text{ wrt } \mathbf{H}) \quad (30)$$

$$\mathbf{M}/\mathbf{E} = \mathbf{H} - \mathbf{G}\mathbf{E}^{-1}\mathbf{F} \quad (31)$$

- Matrix inversion lemma (Sherman-Morrison formula)
- Matrix determinant lemma:

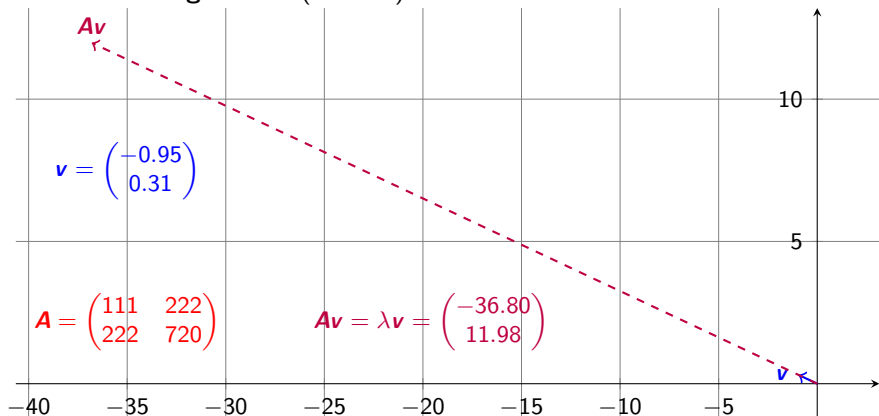
$$|\mathbf{A} + \mathbf{u}\mathbf{v}^\top| = (1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u}) |\mathbf{A}| \quad (32)$$

Eigenvectors and eigenvalues (review)

A vector \mathbf{v} is an eigenvector of a matrix A if its transformation by A scales rather than changes the direction of \mathbf{v} :

$$A\mathbf{v} = \lambda\mathbf{v} \quad (33)$$

where λ is the **eigenvalue** (a scalar).



Eigenvalue-related properties

- The solutions of the characteristic equation

$$\det(\lambda \mathbf{I} - \mathbf{A})\mathbf{v} = \mathbf{0}, \quad \mathbf{v} \neq \mathbf{0} \quad (34)$$

are the eigenvalues λ_i of \mathbf{A} .

- The trace of a matrix is the sum of its eigenvalues

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i \quad (35)$$

- The determinant of a matrix is the product of its eigenvalues

$$\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i \quad (36)$$

- The rank of a matrix is equal to the number of its non-zero eigenvalues

Eigenvalue decomposition (EVD)

We can write:

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{\Lambda} \quad (37)$$

where the columns of $\mathbf{U} \in \mathbb{R}^{n \times n}$ are the eigenvectors of \mathbf{A}

$$\mathbf{U} \in \mathbb{R}^{n \times n} = \begin{pmatrix} | & | & \dots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_n \\ | & | & & | \end{pmatrix} \quad (38)$$

and $\mathbf{\Lambda}$ is a diagonal matrix of the eigenvalues:

$$\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n) \quad (39)$$

- \mathbf{A} is said to be diagonalizable
- EVD: If \mathbf{U} is invertible, then we can also write:

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} \quad (40)$$

Symmetric matrices

When \mathbf{A} is real and symmetric, then \mathbf{U} is orthogonal:

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{\top} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top} \quad (41)$$

And

$$\mathbf{A}^{-1} = \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{U}^{\top} = \sum_{i=1}^d \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^{\top} \quad (42)$$

Singular value decomposition (SVD)

The singular value decomposition of \mathbf{X} generalizes EVD to rectangular matrices:

$$\mathbf{X} = \mathbf{U}\mathbf{\Gamma}\mathbf{V}^T \quad (43)$$

where:

- \mathbf{X} is an $n \times p$ data matrix
- \mathbf{U} is an $n \times p$ orthogonal¹ matrix. The columns of \mathbf{U} are called *left singular vectors*
- $\mathbf{\Gamma}$ is a $p \times p$ diagonal matrix (whose elements are called *singular values*)
- \mathbf{V} is an $p \times p$ orthogonal² matrix. The columns of \mathbf{V} are called *right singular vectors*
- The columns of $\mathbf{U}\mathbf{\Gamma}$ are called the **principal components** of \mathbf{X} .

¹i.e. $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ and $\mathbf{U}^T = \mathbf{U}^{-1}$

²i.e. $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and $\mathbf{V}^T = \mathbf{V}^{-1}$

SVD (cont.)

$$\underbrace{\begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix}}_{\mathbf{X}} = \underbrace{\begin{pmatrix} u_{11} & \cdots & u_{1p} \\ \vdots & \ddots & \vdots \\ u_{n1} & \cdots & u_{np} \end{pmatrix}}_{\mathbf{U}: \text{eigenvectors of } \mathbf{X}\mathbf{X}^T} \underbrace{\begin{pmatrix} \sqrt{\lambda_1} & \cdots & 0 \\ 0 & \ddots & 0 \\ 0 & \cdots & \sqrt{\lambda_p} \end{pmatrix}}_{\substack{\boldsymbol{\Gamma}: \sqrt{\text{eigenvalues of } \mathbf{X}\mathbf{X}^T} \\ \text{also singular values of } \mathbf{X}}} \underbrace{\begin{pmatrix} v_{11} & \cdots & v_{1p} \\ \vdots & \ddots & \vdots \\ v_{p1} & \cdots & v_{pp} \end{pmatrix}}_{\mathbf{V}: \text{eigenvectors of } \mathbf{X}^T\mathbf{X}}$$

- The columns $\mathbf{u}_1, \dots, \mathbf{u}_p$ are the left singular vectors of \mathbf{X}
- The columns $\mathbf{v}_1, \dots, \mathbf{v}_p$ are the right singular vectors of \mathbf{X}
- The elements $\sqrt{\lambda_1} \geq \dots \geq \sqrt{\lambda_p} = 0$ are the singular values of \mathbf{X}
- $\lambda_1 \geq \dots \geq \lambda_p = 0$ are the eigenvalues of $\mathbf{X}\mathbf{X}^T$ and also of $\mathbf{X}^T\mathbf{X}$

Solving a system of equations

A system of linear equations can be represented and solved as follows:

$$\begin{aligned}
 \mathbf{Ax} &= \mathbf{b} \\
 \mathbf{A}^{-1}\mathbf{Ax} &= \mathbf{A}^{-1}\mathbf{b} \\
 \mathbf{I}_n\mathbf{x} &= \mathbf{A}^{-1}\mathbf{b} \\
 \mathbf{x} &= \mathbf{A}^{-1}\mathbf{b}
 \end{aligned}$$

- If $\mathbf{A} \in \mathbb{R}^{m \times n}$, then there are m equations and n unknowns
- If $m = n$, then \mathbf{A} is full rank, and there is a unique solution
- If $m < n$, the system is **underdetermined** (no unique solution)
- If $m > n$, the system is **overdetermined** (no exact solution; solve via least squares)

Least squares

Given a system $\mathbf{Ax} = \mathbf{b}$, the least squares objective is:

$$\min_{\mathbf{x}} f(\mathbf{x}) \equiv \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \quad (44)$$

The gradient is given by:

$$\mathbf{g}(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) = \mathbf{A}^\top \mathbf{Ax} - \mathbf{A}^\top \mathbf{b} \quad (45)$$

The optimal $\hat{\mathbf{x}}$ is found by solving for $\mathbf{g}(\mathbf{x}) = \mathbf{0}$:

$$\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b} \quad (46)$$

And thus,

$$\hat{\mathbf{x}} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \quad (47)$$

is the ordinary least squares (OLS) solution. Checking that the Hessian $\mathbf{H}(\mathbf{x}) = \mathbf{A}^\top \mathbf{A}$ is pd confirms the solution is unique.

Reading assignments

- **PMLI 7**
- **PMLCE 2**