

The Dual Perspective: The Hidden Geometry of $y = Ax$



Tomio Kobayashi

Following

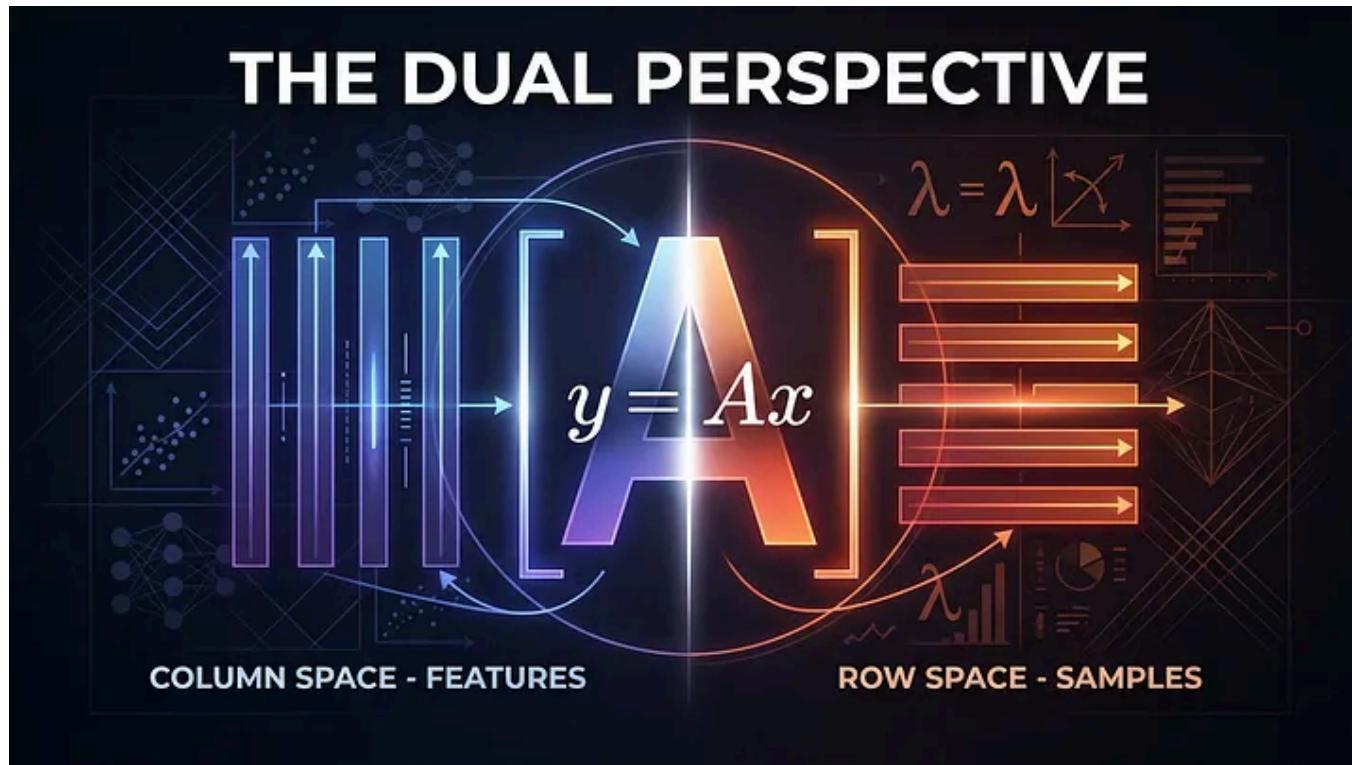
14 min read · 2 hours ago



10



...



The Fabric of Data

In the vast landscape of quantitative science — spanning physics, econometrics, machine learning, and signal processing — there is one mathematical structure that stands above all others in its ubiquity: the linear system.

$$\mathbf{y} = A\mathbf{x}$$

At first glance, this equation appears deceptively simple. We are taught in introductory courses to view the matrix A as a static container — a rectangular grid of numbers that functions as a mechanical operator, transforming an input vector x into an output vector y . We memorize rules for multiplication and inversion, treating the matrix as a monolithic block.

However, to truly understand the mechanics of modern data science and the statistical tools that drive it, we must abandon this static view. We must theorize that a matrix is not a single object, but rather a superposition of two distinct, overlapping realities. It is, simultaneously, a collection of **Row Vectors** and a collection of **Column Vectors**.

These two sets are not disjoint. They occupy the exact same coordinate space and describe the exact same data elements. They are the “warp and weft” of the data fabric. Yet, depending on which thread we pull — the row or the column — we extract completely different information. One perspective reveals the statistical variance of our features; the other reveals the geometric similarity of our samples.

The fundamental “problem” in data analysis — and the source of its richness — is that we are almost always dealing with matrices where these two perspectives are unbalanced. We rarely encounter the perfect square matrices of textbook algebra. We deal with “tall” matrices (more data than variables) or “fat” matrices (more variables than data).

This article explores how adopting this **Dual Perspective** allows us to navigate these shapes. We will see how deriving the symmetric interactions of these views — $A^T A$ and $A A^T$ — provides the natural statistical engine for the canonical toolsets of the field: from Linear Regression and the Pseudo-Inverse to PCA, SVD, and Kernel Methods.

Part 1: The Anatomy of $y=Ax$

To understand the duality, we must first break down the fundamental equation $y=Ax$ into its two component interpretations. Let us assume A is a matrix with m rows and n columns.

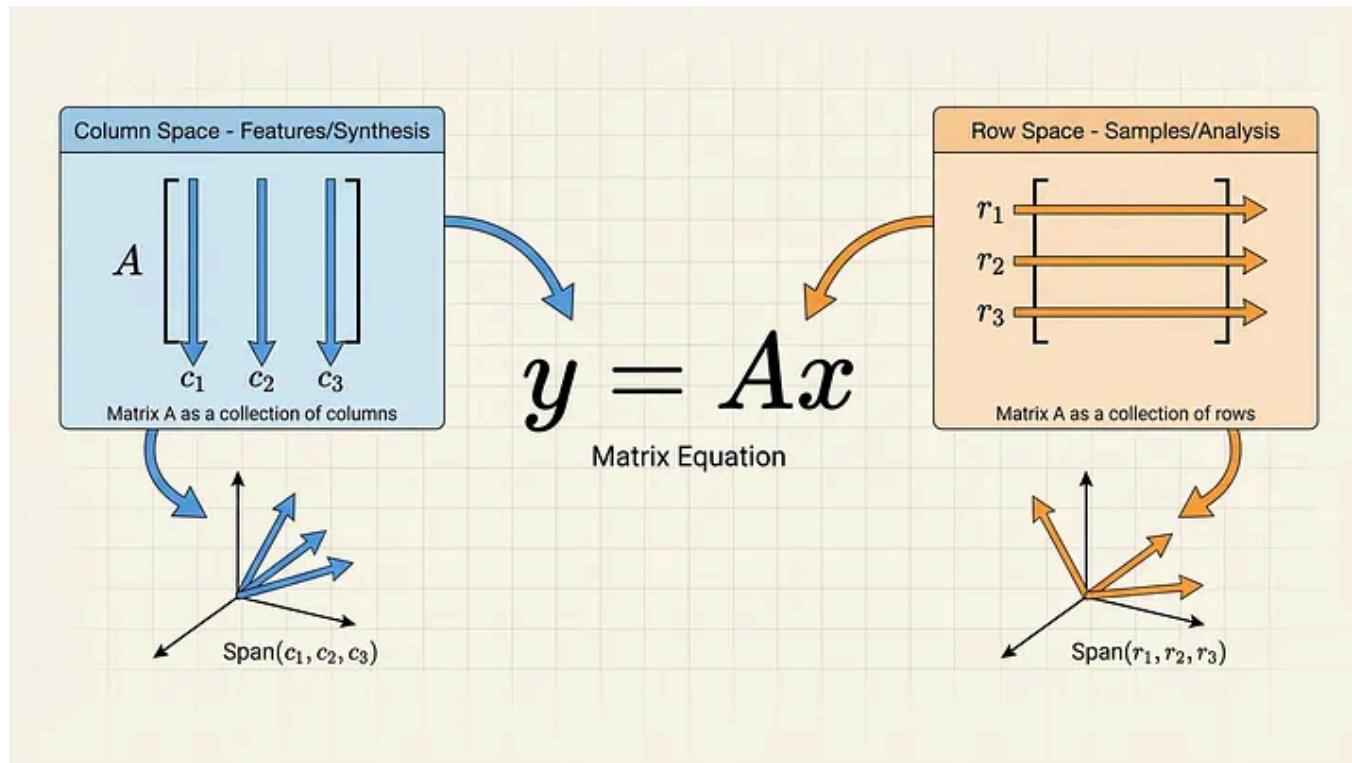


Figure 1: The dual perspective of matrix A – simultaneously viewed as column vectors (features/synthesis) and row vectors (samples/analysis)

Perspective A: The Column Space (The “What”)

In data science, the columns of a matrix usually represent the **features**, dimensions, or sensors. For example, in a housing dataset, columns might represent “Square Footage,” “Number of Bedrooms,” and “Year Built.”

When we view A as a set of column vectors $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n$, the operation $\mathbf{y} = A\mathbf{x}$ takes on a specific meaning. It is a **Linear Combination**.

$$\mathbf{y} = x_1\mathbf{c}_1 + x_2\mathbf{c}_2 + \cdots + x_n\mathbf{c}_n$$

Here, the unknown vector \mathbf{x} acts as a set of controls or weights. We are asking: *How much of Feature 1 plus how much of Feature 2 do we need to mix together to reproduce the result \mathbf{y} ?*

This is the perspective of **Synthesis**. We are trying to build the output \mathbf{y} using the building blocks provided by the columns of A . This view dominates fields like regression and signal reconstruction.

Perspective B: The Row Space (The “Who”)

Simultaneously, the rows of A usually represent the **samples**, observations, or data points. In our housing example, each row is a specific house.

When we view A as a set of row vectors $\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_m$, the operation $\mathbf{y} = A\mathbf{x}$ changes meaning. It becomes a batch of **Dot Products**.

$$y_i = \mathbf{r}_i \cdot \mathbf{x}$$

Here, the calculation represents a projection or a similarity test. We are asking: *How well does the input vector \mathbf{x} align with Sample 1? How well does it align with Sample 2?*

This is the perspective of **Analysis**. We are scanning the database of samples against a probe vector \mathbf{x} to see how they interact. This view dominates fields like search engines, pattern recognition, and classification.

The Shape Dilemma

The friction in linear algebra arises because m and n are rarely equal.

MATRIX SHAPES IN LINEAR ALGEBRA

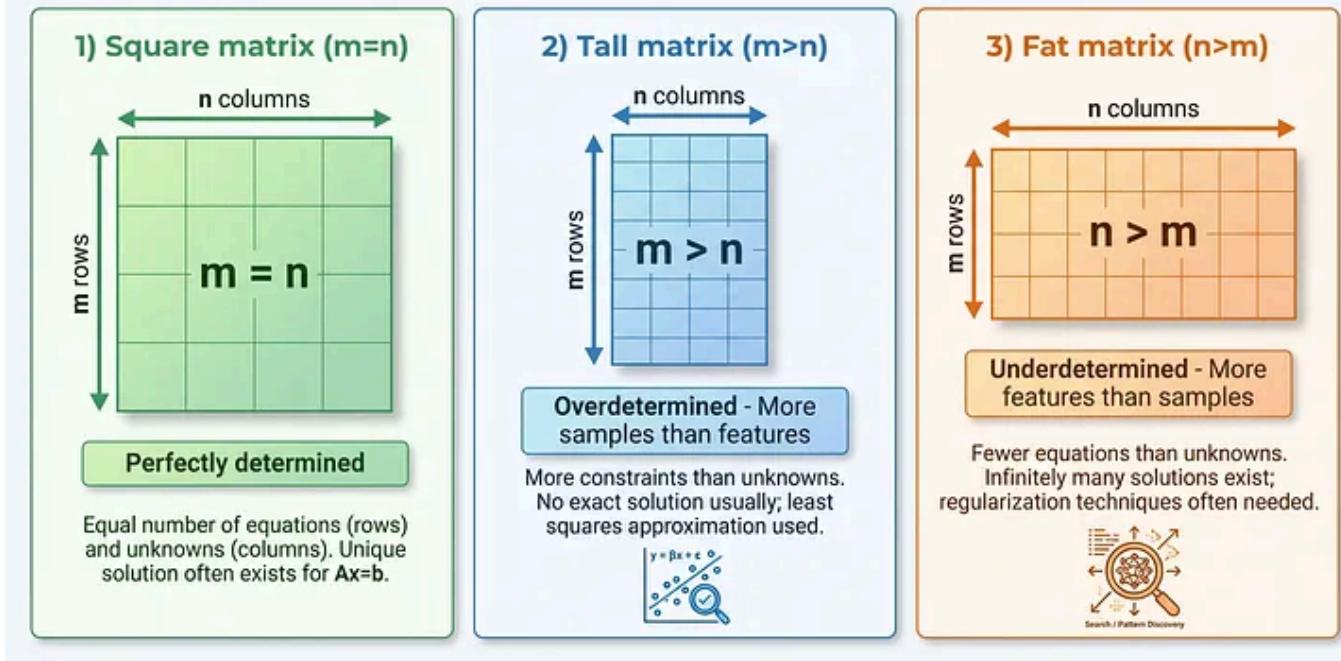


Figure 2: The three matrix configurations – square (perfectly determined), tall (overdetermined), and fat (underdetermined)

1. **The Overdetermined System ($m>n$):** We have 1,000 houses (rows) but only 3 features (columns). In the column perspective, we are trying to reach a vector y in a 1,000-dimensional space using only 3 building blocks. It is geometrically impossible to hit the target exactly; y lies outside the “span” of the columns.

2. The Underdetermined System ($n > m$): We have 20 houses (rows) but 50,000 genetic markers (columns). In the column perspective, we have too many building blocks. We can combine them in infinite ways to reproduce y .

To solve these problems, we cannot look at the rectangular matrix A alone. We must look at how the rows interact with the rows, and how the columns interact with the columns.

Part 2: Deriving the Interaction Matrices

The bridge between raw data and statistical insight is the **Dot Product**. The dot product is the atom of linear algebra — it reduces two vectors into a single number that describes their relationship.

By applying the dot product systematically to our dual perspectives, we derive two new matrices: $A^T A$ and $A A^T$. These are not merely algebraic rearranging; they are the transition from “Data” to “Statistics.”

The Column Interaction: $A^T A$ (The Covariance Matrix)

If we multiply the transpose of A by A ($A^T A$), we are mathematically dotting every column against every other column. The result is an $n \times n$ square, symmetric matrix.

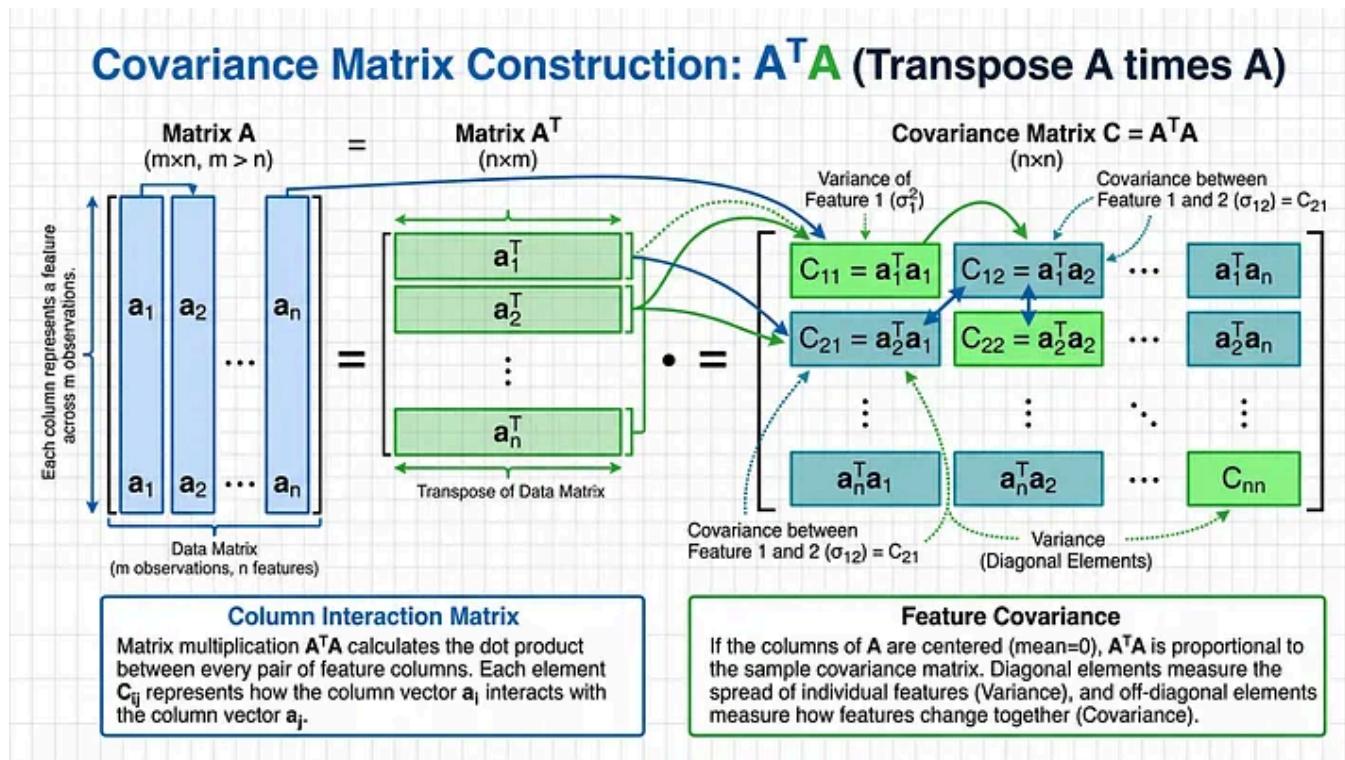


Figure 3: The formation of $A^T A$ – the column interaction matrix revealing feature covariance

If we assume our data is centered (mean subtracted), $A^T A$ is proportional to the Covariance Matrix.

1. The Diagonal (Variance): The element at (i,i) is the dot product of column i with itself: $\mathbf{c}_i \cdot \mathbf{c}_i = \| \mathbf{c}_i \|^2$.

- **Statistical Meaning:** This is the **Variance** of feature i . It represents the “energy” or “information content” of that specific variable. In signal processing, this is the power of the signal.

2. The Off-Diagonal (Covariance): The element at (i,j) is $\mathbf{c}_i \cdot \mathbf{c}_j$.

- **Statistical Meaning:** This is the **Covariance** between feature i and feature j . It tells us if these two variables move together. If the value is zero, the features are **Orthogonal** (uncorrelated).

The Takeaway: $A^T A$ maps the internal structure of our variables. It tells us which features are noisy (high variance) and which are redundant (high covariance).

The Row Interaction: AA^T (The Gram Matrix)

If we multiply A by its transpose (AA^T), we are dotting every row against every other row. The result is an $m \times m$ square, symmetric matrix.

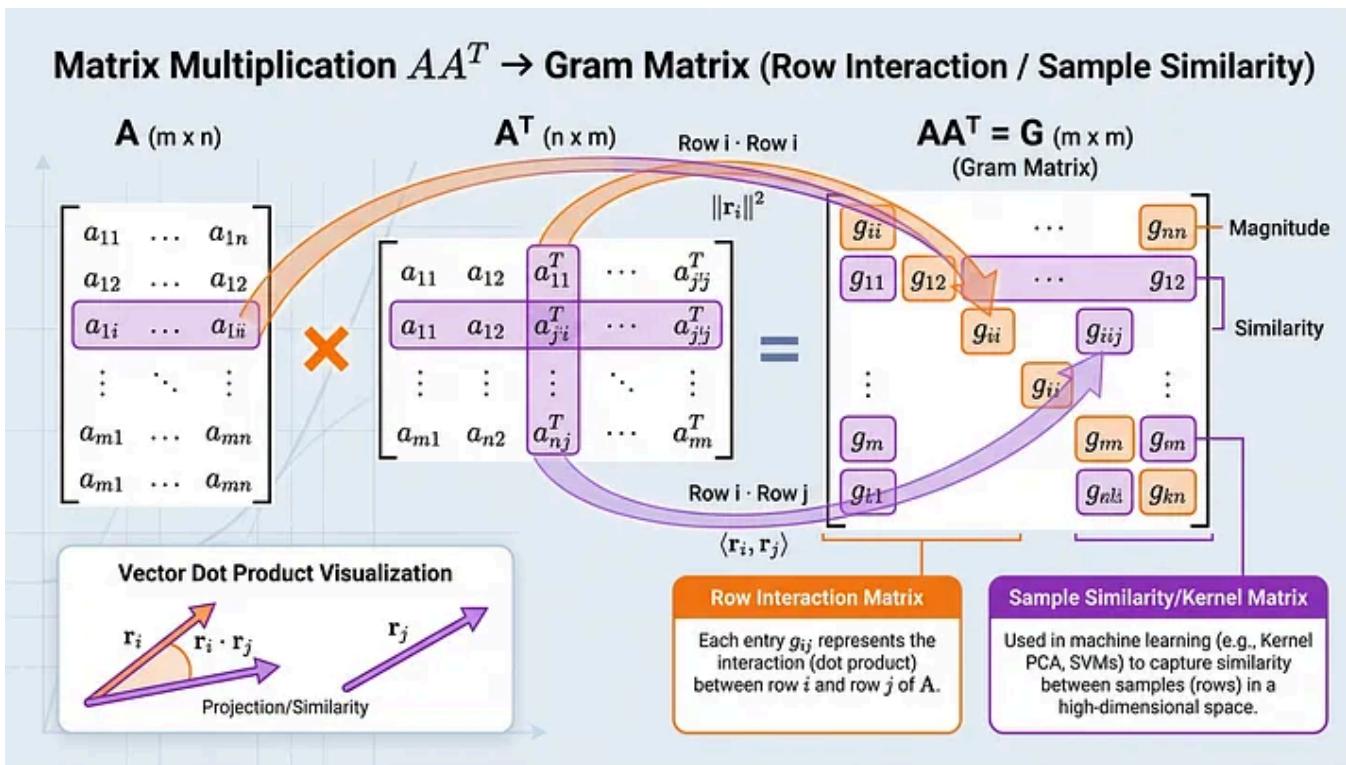


Figure 4: The formation of AA^T – the row interaction matrix revealing sample similarity

This is often called the **Gram Matrix** or the **Kernel Matrix**.

1. **The Diagonal (Magnitude):** The element at (i,i) is $r_i \cdot r_i$.

- **Statistical Meaning:** This is the squared magnitude of sample i . It tells us the “intensity” of that specific observation.

1. The Off-Diagonal (Similarity): The element at (i,j) is $\mathbf{r}_i \cdot \mathbf{r}_j$.

- **Statistical Meaning:** This is the geometric **Similarity** between sample i and sample j . If the vectors are normalized, this is exactly the **Cosine Similarity**.

The Takeaway: AA^T maps the geometric relationship between our subjects. It tells us which samples are clustered together and which are outliers.

The Statistical Engine: Practical Interpretations

Before diving into advanced algorithms, we must recognize that A^TA and AA^T are not abstract mathematical constructs — they are the **mechanical calculators** of fundamental statistics. In data engineering and signal processing, these matrices provide direct, computable measurements that answer critical questions about our data.

1. Variance and Covariance as Signal Quality Metrics

When we compute A^TA for centered data, we obtain:

$$\Sigma = \frac{1}{m-1} A^T A$$

This covariance matrix serves as a “**Control Panel**” for understanding feature quality:

The Diagonal Elements (Variance = Signal Energy):

- In engineering terms, $c_i \cdot c_i$ measures the **energy** of the signal from sensor i
- In statistical terms, this is the **variance** — how much information or spread exists in feature i
- **Practical Question Answered:** “Which sensors are providing strong signals vs. which are just noise?”

The Off-Diagonal Elements (Covariance = Signal Redundancy):

- The value $c_i \cdot c_j$ measures how two features move together
- **Positive covariance:** Features increase together (redundant information)
- **Zero covariance:** Features are **orthogonal** — the “Holy Grail” of independent sensors
- **Practical Question Answered:** “Are my sensors measuring the same underlying phenomenon, or providing unique information?”

Engineering Example: In a vibration monitoring system with multiple accelerometers, $A^T A$ tells us which sensors are providing unique data and which are redundantly measuring the same vibration mode.

PCA Connection: Principal Component Analysis is simply the algorithmic process of rotating our coordinate system so that all the energy (variance) aligns with the diagonal axes and all redundancy (covariance) becomes zero. We maximize information and eliminate redundancy.

2. Similarity and Cosine Distance as Pattern Recognition

When we compute $A^T A$, we transition from statistics to **geometry** and **pattern matching**:

The Dot Product as Similarity Measure:

- Each entry (i,j) computes $\mathbf{r}_i \cdot \mathbf{r}_j$ – the alignment between samples
- **High value:** Samples i and j have significant overlap in their feature patterns
- **Low/zero value:** Samples are unrelated or orthogonal in feature space
- **Practical Question Answered:** “Which observations are similar to each other?”

Normalized Similarity (Cosine Similarity): If we normalize each row to unit length before multiplication, AA^T strictly contains cosine similarities:

$$\text{Cosine Similarity} = \cos(\theta) = \frac{\mathbf{r}_i \cdot \mathbf{r}_j}{\|\mathbf{r}_i\| \cdot \|\mathbf{r}_j\|}$$

- Value of **1.0**: Samples are identical in direction
- Value of **0.0**: Samples are completely independent
- Value of **-1.0**: Samples are opposite

Recommendation System Example: In Netflix's collaborative filtering, rows represent users and columns represent movie ratings. Computing AA^T reveals which users have similar taste patterns. If users i and j have high dot product, they rated similar movies highly — the system can recommend movies user i liked to user j .

Search Engine Example: In document retrieval, rows are documents and columns are word frequencies. When you submit a query (a vector), the search engine computes the dot product between your query vector and each document row — effectively computing a slice of AA^T . Documents with the

highest dot product (smallest angle, highest cosine similarity) are the best matches.

3. The Pseudo-Inverse as Statistical Normalization

The least squares solution to $y = Ax$ demonstrates how these interaction matrices combine:

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{y}$$

Breaking down the components reveals the statistical logic:

1. $A^T \mathbf{y}$: Computes the **correlation** between each feature and the target

- Measures raw similarity: “How much does each feature look like the desired output?”

2. $(A^T A)^{-1}$: Inverts the **covariance matrix**

- Acts as a normalization filter: “Divide by variance to remove the effect of feature scale”

- Accounts for redundancy: “Adjust for the fact that features might be correlated”

Engineering Translation: The solution is fundamentally:

$$\mathbf{x} = \frac{\text{Correlation}}{\text{Variance}}$$

We take the raw similarity between inputs and outputs, then normalize by how “spread out” or “energetic” the inputs are. The pseudo-inverse is a **statistical filter** that extracts signal from noise by leveraging the dual perspective — it uses column statistics ($A^T A$) to optimally weight the observed correlations ($A^T y$).

Signal Processing Example: When fitting a regression line through noisy sensor measurements, we’re not just finding “any” line — we’re finding the line that maximally explains the variance in our output while accounting for the variance and covariance structure of our input sensors.

Summary of the Statistical Engine

By recognizing the dual perspective, the matrices $A^T A$ and $A A^T$ reveal their true nature:

- **$A^T A$ (Column Perspective): The Variance/Covariance Machine**
 - Measures signal quality and feature independence
 - Answers: “Which sensors are informative?” and “Which sensors are redundant?”
 - Applications: PCA, feature selection, signal processing, dimensionality reduction
- **$A A^T$ (Row Perspective): The Similarity/Correlation Machine**
 - Measures pattern matching and geometric relationships
 - Answers: “Which observations are similar?” and “How do samples cluster?”
 - Applications: Recommendation systems, search engines, clustering, nearest neighbor

Data analysis is the practice of intelligently toggling between these two views – extracting feature statistics when we need to understand variable quality, and computing sample similarities when we need to find patterns and make predictions. The power of linear algebra lies not in memorizing formulas,

but in understanding that these interaction matrices are the computational engines that transform raw data into actionable statistical insight.

Part 3: The Unifying Theory: Spectral Invariance

Before diving into the toolsets, we must establish why we can freely switch between these two views. It stems from a profound theorem in linear algebra regarding Eigenvalues.

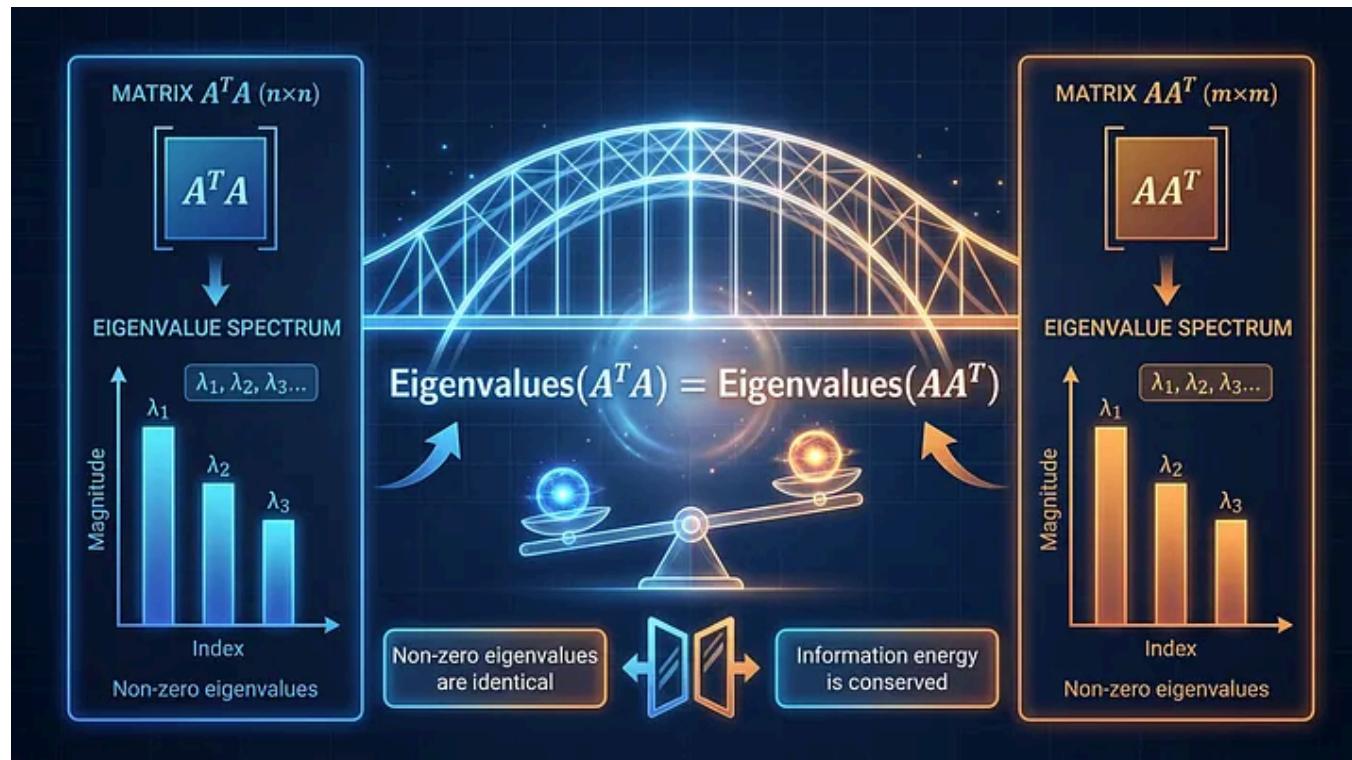


Figure 5: Spectral invariance — the non-zero eigenvalues of $A^T A$ and AA^T are identical, demonstrating information conservation

Even though $A^T A$ is an $n \times n$ matrix (describing features) and AA^T is an $m \times m$ matrix (describing samples), they share the exact same set of non-zero eigenvalues.

$$\text{Eigenvalues}(A^T A) = \text{Eigenvalues}(AA^T)$$

The Physical Interpretation: Eigenvalues represent the magnitude of variance or “stretch” along the principal axes of the data. This theorem states that the total amount of “information energy” in the dataset is conserved. It does not matter if you analyze the variance of the columns or the clustering of the rows — the fundamental spectral fingerprint of the data is identical.

This **Spectral Invariance** is the license that allows algorithms like SVD and the Pseudo-Inverse to work. It guarantees that the geometry of the row space and the geometry of the column space are mirror images of each other.

Part 4: The Statistical Toolsets

Now we can see how the major algorithms of data analysis are simply different ways of leveraging these derived interaction matrices.

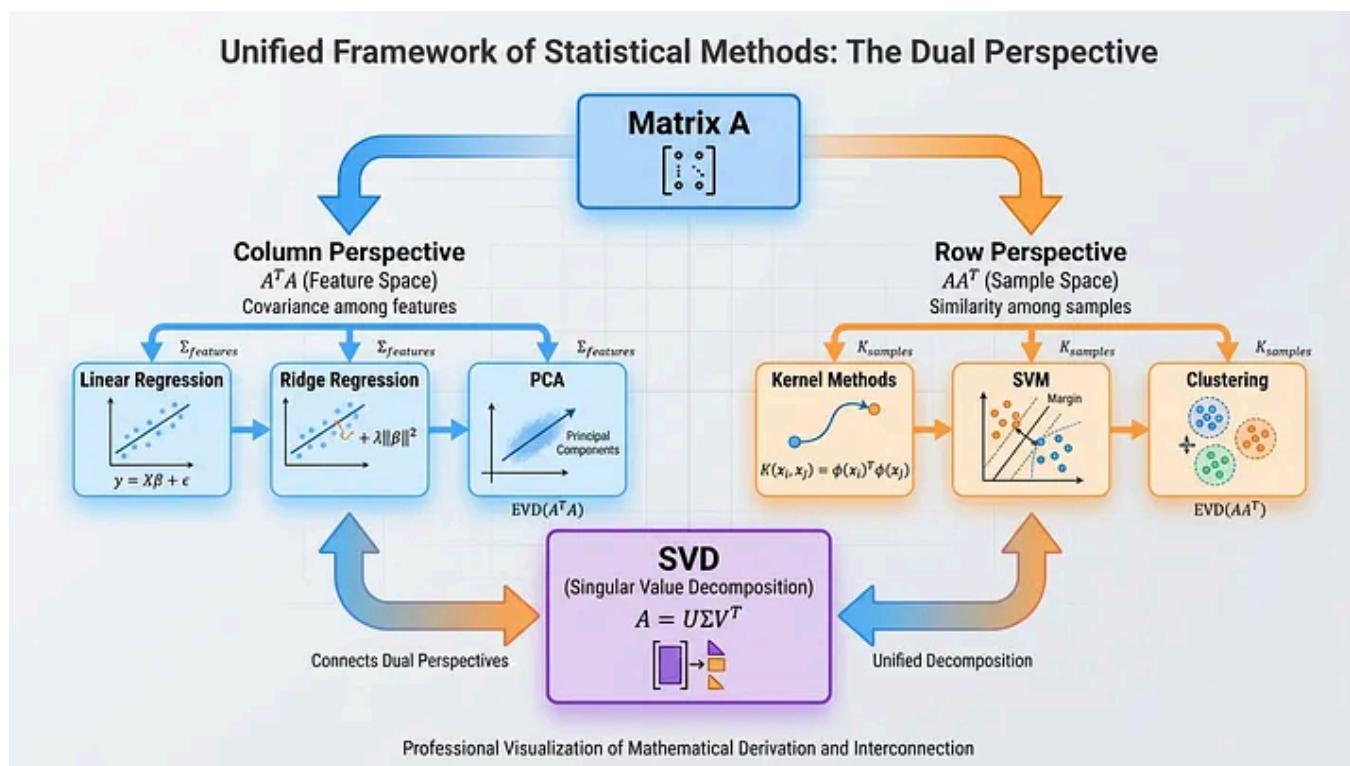


Figure 6: The unified framework showing how statistical methods derive from the dual perspective

1. The Normal Equation (Linear Regression)

Context: The matrix is Tall ($m > n$). We have more equations than unknowns.

The Problem: $y = Ax$ has no solution. The vector y is floating in m -dimensional space, but the columns of A only cover a small n -dimensional subspace.

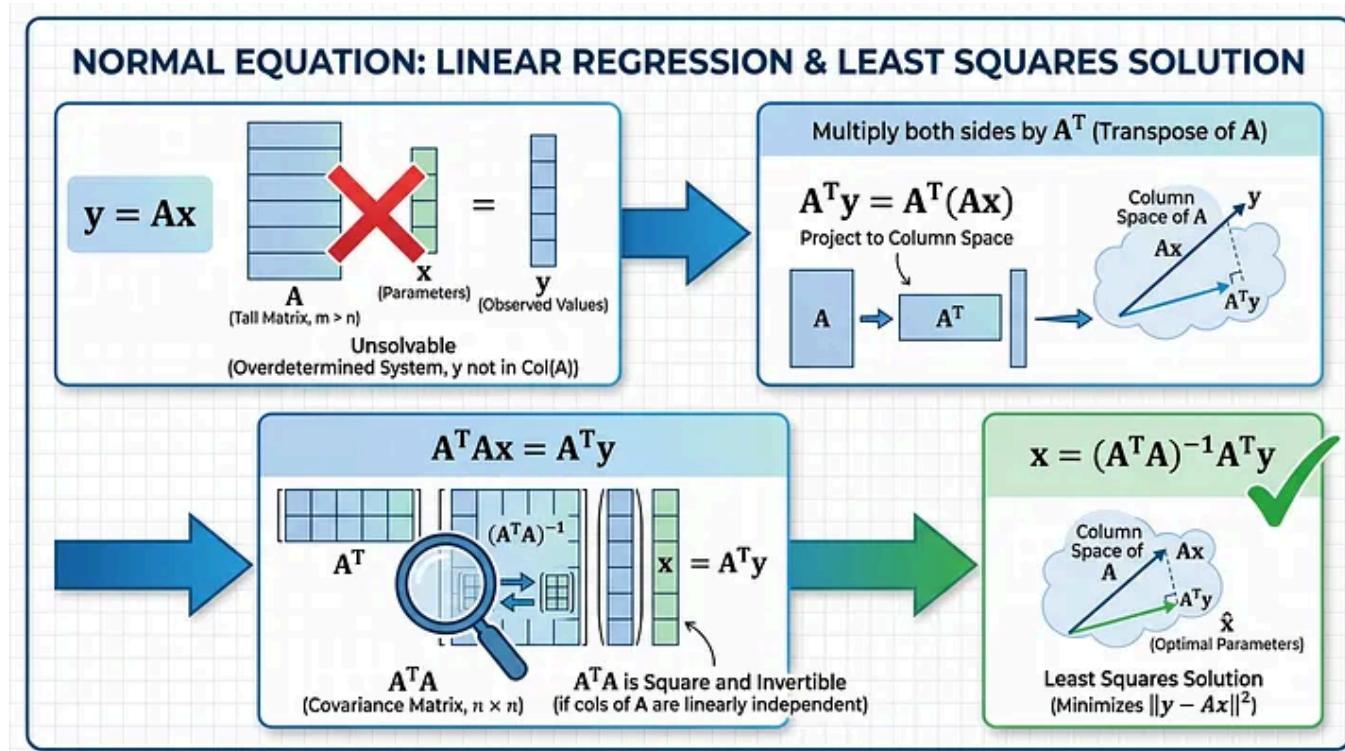


Figure 7: The Normal Equation solution process — projecting onto the column space for least squares regression

The Solution: We project the problem into the **Column Interaction** space. We multiply both sides by A^T .

$$A^T A \mathbf{x} = A^T \mathbf{y}$$

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{y}$$

The Dual Perspective Analysis: Look at the components of this solution:

- $A^T \mathbf{y}$: This vector represents the **Covariance between the Features and the Target**. It asks: “How much does Feature 1 look like the output \mathbf{y} ? ”
- $(A^T A)^{-1}$: This is the **Inverse Covariance Matrix**. It acts as a “whitening” filter. It corrects for the fact that the features might be correlated with each other.

Conclusion: Linear regression effectively calculates the raw correlation between inputs and outputs, and then divides it by the variance of the inputs to isolate the true relationship.

2. Ridge Regression (Regularization)

Context: Sometimes, the columns of A are almost identical

(Multicollinearity). **The Problem:** If two features are highly correlated, the covariance matrix ATA becomes “singular” or “degenerate.” Its determinant approaches zero, and we cannot invert it. The solution becomes unstable; small noise in y leads to massive swings in x .

The Solution: We hack the interaction matrix.

$$\mathbf{x} = (A^T A + \lambda I)^{-1} A^T \mathbf{y}$$

The Dual Perspective Analysis: We add λ to the diagonal of $A^T A$.

- Since the diagonal represents Variance, we are artificially inflating the variance of every feature.
- We are telling the math: “Pretend these features are more distinct and independent than they actually are.”
- This stabilizes the inversion, trading a little bit of accuracy (Bias) for a massive reduction in instability (Variance).

3. Principal Component Analysis (PCA)

Context: We have a dataset with many noisy, correlated variables. We want to simplify it. **The Problem:** The “Column Perspective” is messy. The variances are spread out, and the covariances (off-diagonals) are non-zero.

The Solution: We perform Eigendecomposition on the **Column Interaction** matrix ($A^T A$).

The Dual Perspective Analysis: PCA searches for a new set of axes (eigenvectors) such that if we rotate the data to align with them:

1. All the **Variance** is maximized along the diagonals.
2. All the **Covariance** (off-diagonals) becomes zero.

PCA effectively asks: *What is the ‘true’ coordinate system of this data where every feature is perfectly independent?* It answers this by diagonalizing $A^T A A$.

4. The Moore-Penrose Pseudo-Inverse (A^+)

Context: We want a universal way to solve $y = Ax$ for any matrix shape (Tall, Fat, or Square).

The Solution: The Pseudo-Inverse (A^+) is the generalized solver. Crucially, its calculation depends entirely on which perspective is dominant.

- **Case 1: Tall Matrix ($m > n$).** The constraint is in the Feature Space.

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

We invert the **Feature Covariance**. This is the Least Squares solution.

- **Case 2: Fat Matrix ($n > m$).** The constraint is in the Sample Space.

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

We invert the **Sample Similarity**. This is the Minimum Norm solution.

The Dual Perspective Analysis: The Pseudo-Inverse proves that you cannot solve a rectangular problem without anchoring yourself in one of the two symmetric squares. You effectively “project” the problem into whichever symmetric domain ($n \times n$ or $m \times m$) is smaller and invertible, solve it there, and map it back.

5. Kernel Methods and SVMs

Context: We have a “Fat” matrix or we want to project data into infinite dimensions to make it separable (like in a Support Vector Machine).

The Problem: We cannot calculate the columns in an infinite-dimensional space. Computing $A^T A$ is impossible.

The Solution: We rely exclusively on the **Row Interaction** matrix (AA^T).

The Dual Perspective Analysis: The “Kernel Trick” is the realization that many algorithms (like SVMs) *only* involve the dot products between samples (AA^T). They never actually need the individual columns of AA . We replace the computed matrix AA^T with a function $K(r_i, r_j)$ that computes the similarity directly.

- Instead of saying “Here are the features of Sample A and Sample B,” we say “Here is how similar Sample A is to Sample B.”
- By working entirely in the Row/Similarity perspective, we bypass the need to define the columns explicitly, allowing for powerful non-linear classification.

6. Singular Value Decomposition (SVD)

The Solution: SVD breaks A into three matrices: $A=U\Sigma V^T$.

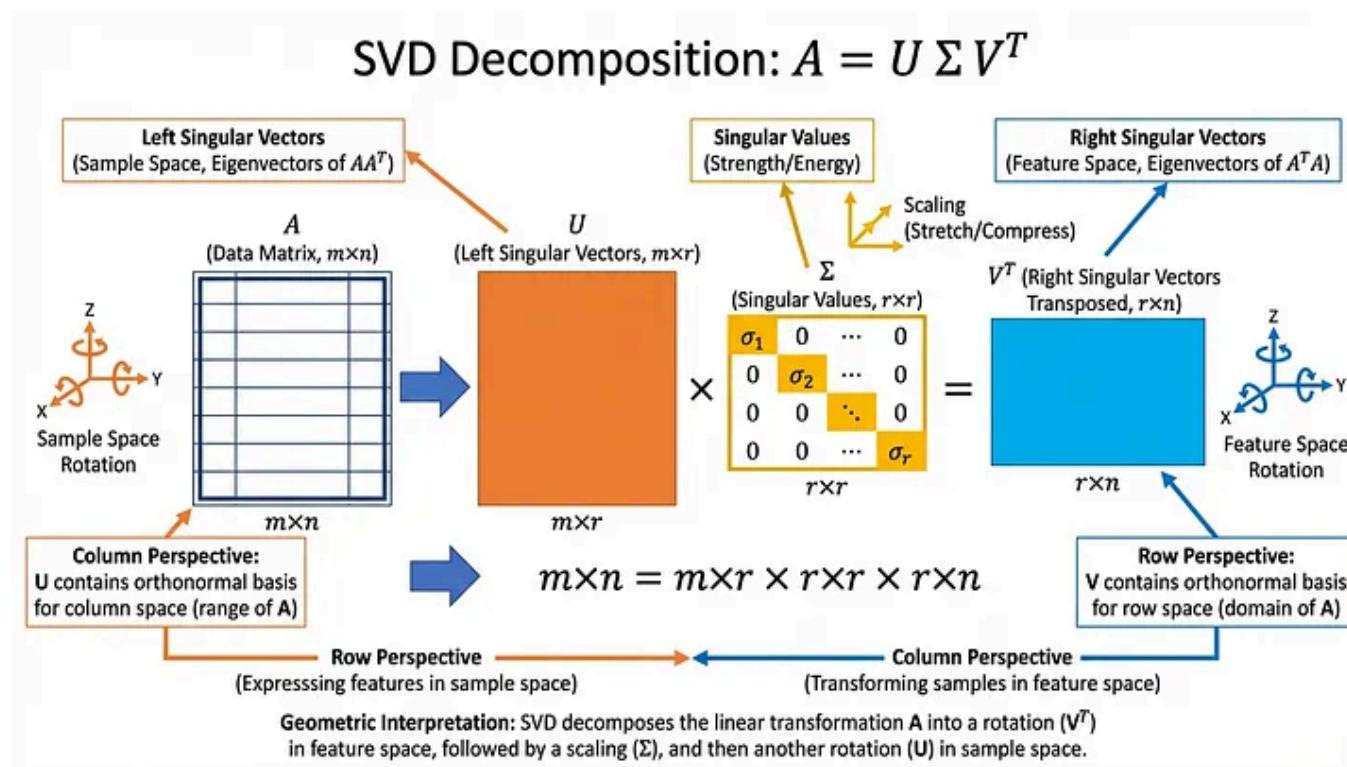


Figure 8: Singular Value Decomposition — simultaneously solving both row and column perspectives

The Dual Perspective Analysis: SVD is the algorithm that solves the Row Perspective and Column Perspective simultaneously and links them

together.

- **V (Right Singular Vectors):** These are the Eigenvectors of $A^T A$. They describe the principal axes of the **Features**.
- **U (Left Singular Vectors):** These are the Eigenvectors of $A A^T$. They describe the principal axes of the **Samples**.
- **Σ (Singular Values):** These are the square roots of the shared eigenvalues. They represent the **Gain** or strength of the connection between the Row patterns and the Column patterns.

SVD allows us to compress a matrix by keeping only the top k singular values. This is equivalent to saying: “Keep only the strongest patterns that exist in *both* the row correlations and the column correlations.”

Conclusion: The Overlapping Reality

The power of linear algebra in data science does not come from memorizing algorithms. It comes from intuition — specifically, the intuition of the **Dual Perspective**.

When we look at a dataset $y = Ax$, we are not looking at a static grid. We are looking at a dynamic interplay between features and samples.

- We use the **Column Perspective** ($A^T A$) to understand the **Variance** and structure of our variables.
- We use the **Row Perspective** ($A A^T$) to understand the **Similarity** and geometry of our observations.

From the basic regression line to the most complex deep learning kernels (which perform operations on massive tensors), the math remains consistent. We are constantly deriving these interaction matrices, analyzing their eigenvalues, and using them to extract signal from noise. By acknowledging that these two perspectives overlap and define each other, we unlock the true analytical power of the matrix.

Mathematics

Software Engineering

Computer Science

Statistics

Machine Learning



Written by Tomio Kobayashi

501 followers · 2 following

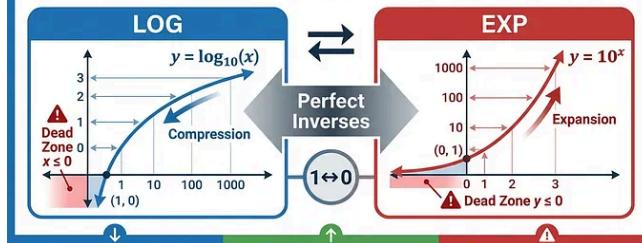
Following ▾

System Architect | Technology Researcher | Master in Computer Science from Georgia Tech | <https://www.linkedin.com/in/tomio-kobayashi-9869ba30/>

More from Tomio Kobayashi

Logarithms & Exponentials

Dead Zones, Fractional Values, and the Mathematics of Scale



Tomio Kobayashi

Understanding Logarithms and Exponentials: The Mathematics o...

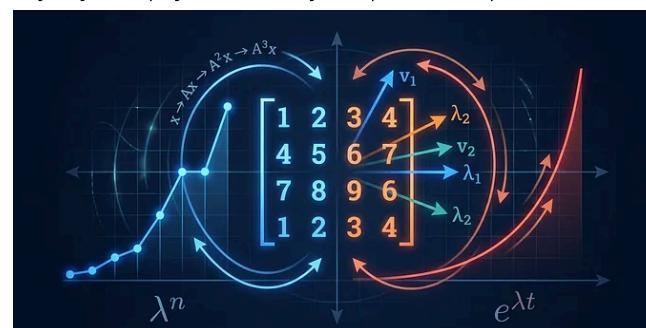
Dead Zones, Fractional Values, and the Mathematics of Scale

Jan 26

200



...



Tomio Kobayashi

Square Matrices as Recurrent Transformation Engines: The...

Matrices as Transformation Engines

4d ago

161



...



Tomio Kobayashi

Exploring Partial Derivatives in Two-Variable Systems: A Fresh...



Tomio Kobayashi

What is this evil “-1” in Sample Variance?

Nov 1, 2025

98



...

Why Sample Variance Divides by $n-1$ Instead of n

Dec 29, 2025

262



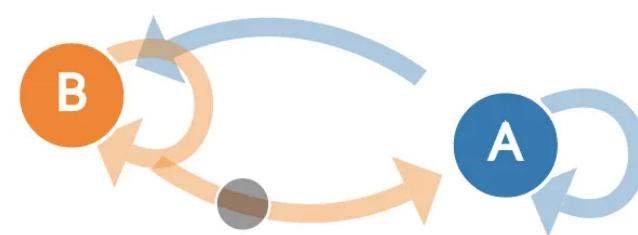
...

See all from Tomio Kobayashi

Recommended from Medium



Basak Kaya

Two Ways of Thinking About Uncertaintyhttps://medium.com/@tomkob99_89317/the-dual-perspective-the-hidden-geometry-of-y-ax-e1e076a528e9

Chirag Dubey

Markov Chains: The Hidden Logic Behind Randomness

A Gentle Introduction to Probability and Statistical Thinking

Jan 10 👏 3 💬 2



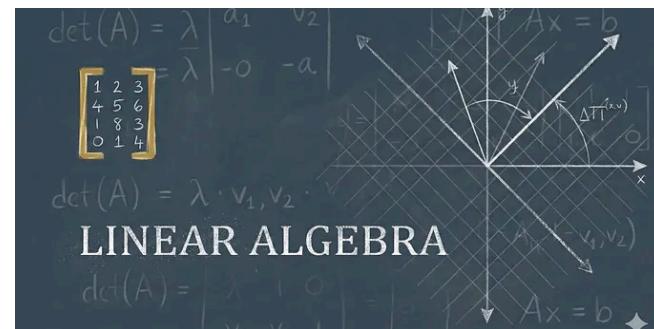
...

Imagine you are standing in a maze. The path is not clear except for what's right in front of...

Jan 24 👏 4



...



Infinity Slap

LA1: Geometry of Linear Algebra

The main idea of any linear algebra course that you take up is to solve a system of linea...

Jan 12 👏 6



...

Tomio Kobayashi

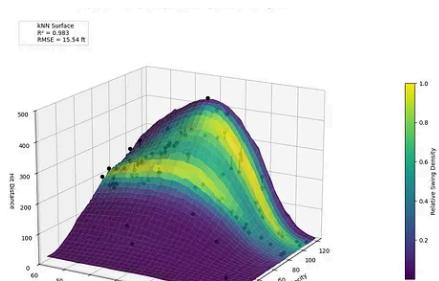
Square Matrices as Recurrent Transformation Engines: The...

Matrices as Transformation Engines

4d ago 👏 161



...





Robbie Dudzinski

When Correlation Fails: Modeling Hit Distance in Baseball

Exit Velocity and Launch Angle have no correlation—but together they predict the...

6d ago

69

1



...



Kiran vutukuri

81. Large Language Models: The Illusion of Understanding

LLM stands for Large Language Model. It's a type of artificial intelligence trained on...



Jan 15

50



...

[See more recommendations](#)