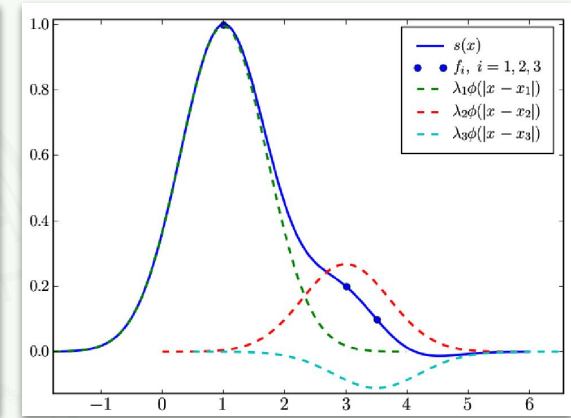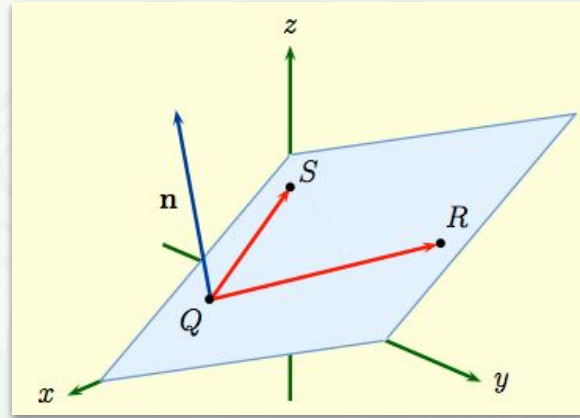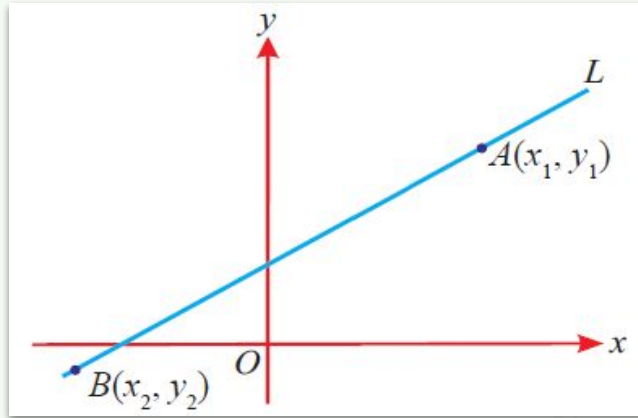# Linear Algebra II

**Michael S. Murillo**
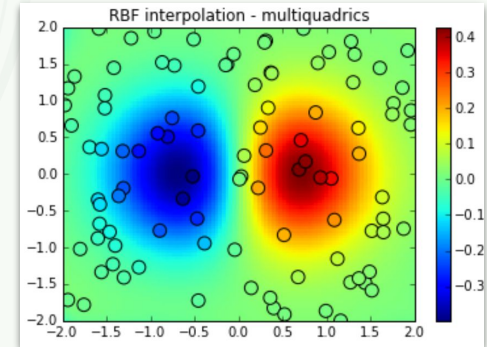
Computational Mathematics, Science and Engineering
Michigan State University

CMSE

# Review of Linear Regression



$$\mathbf{y} = K\mathbf{w} \qquad \mathbf{w} = K^{-1}\mathbf{y}$$

CMSE

# Multiple Linear Regression

In data science, we usually have many input features (independent variables):

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \ldots$$

In general, every additional feature you add gives you more predictive power. There are two caveats:

1. The feature $x_n$ is irrelevant. In this case, $w_n = 0$; little harm done.

2. If two features are linearly dependent, you need to drop one of them. It adds no new information and it will likely cause mathematical problems.
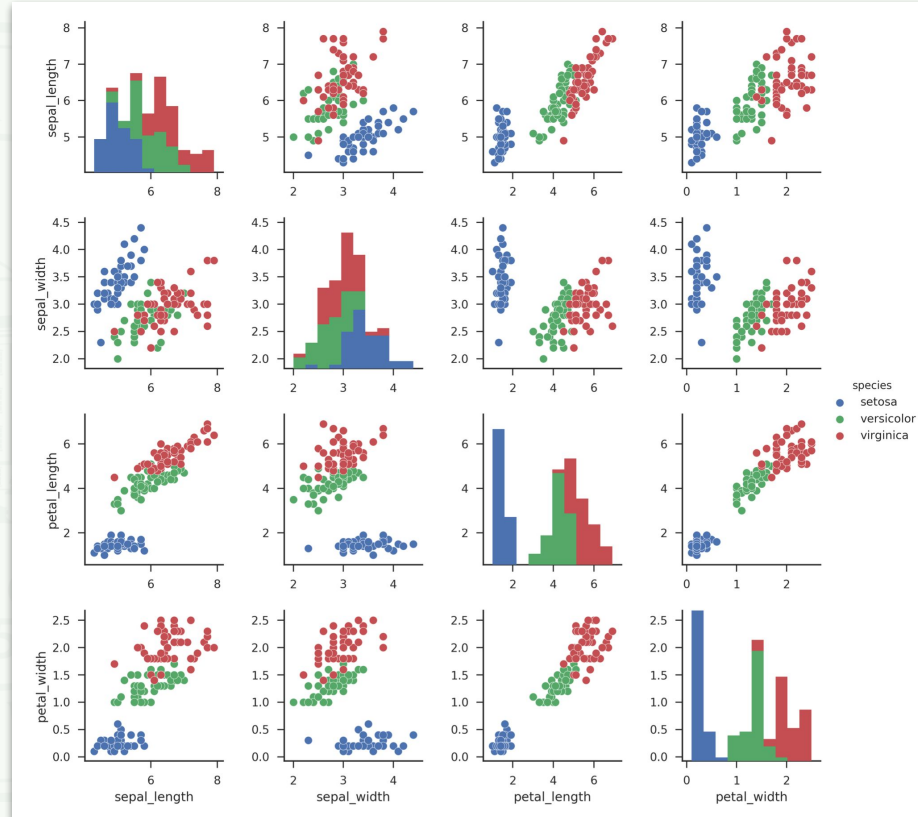
CMSE

# Multiple Linear Regression: Iris

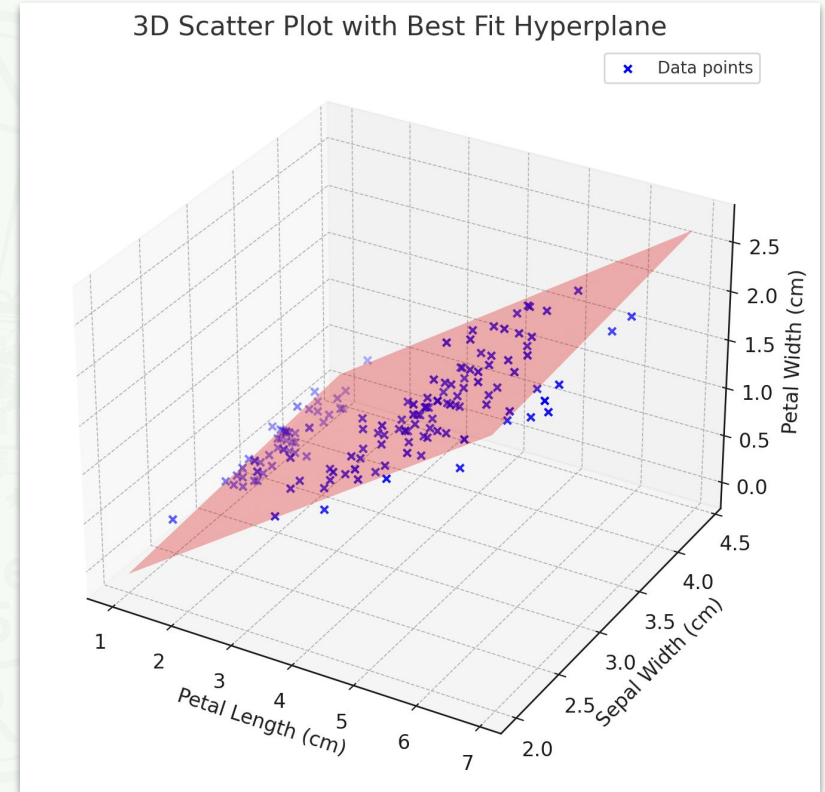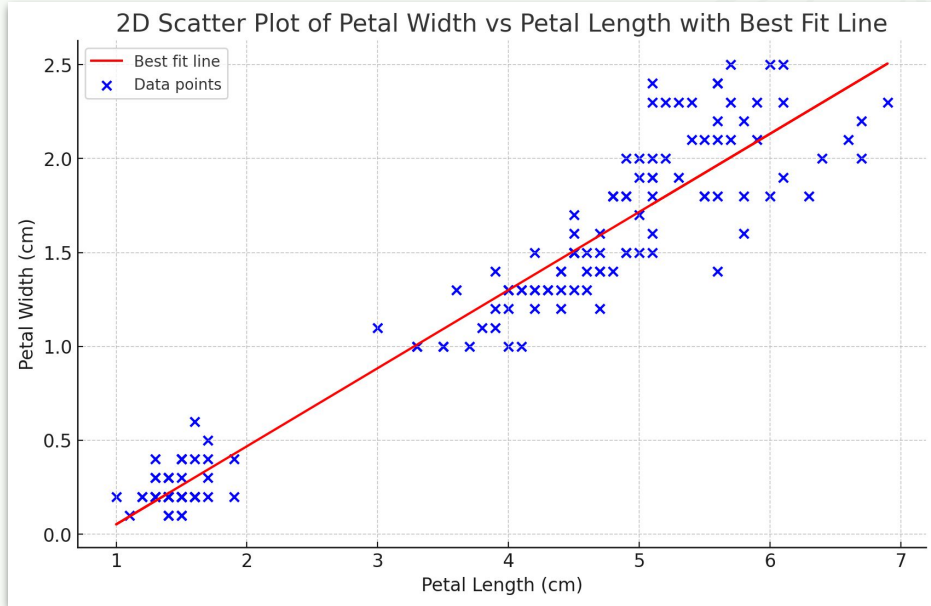Suppose we want to predict petal width for versicolor (green).

petal_width = $w_0$

petal_width = $w_0 + w_1$ petal_length

petal_width = $w_0 + w_1$ petal_length + $w_2$ sepal_width
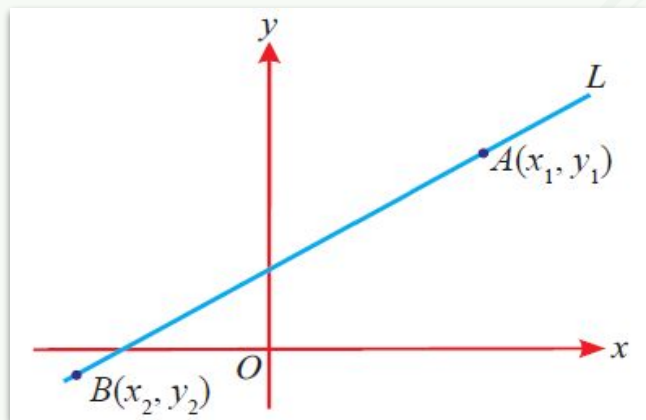
Note that petal_width versus sepal_length is fairly flat.

CMSE

# Multiple Linear Regression with Iris



2D Scatter Plot of Petal Width vs Petal Length with Best Fit Line



3D Scatter Plot with Best Fit Hyperplane

CMSE

# Linear Algebra Details

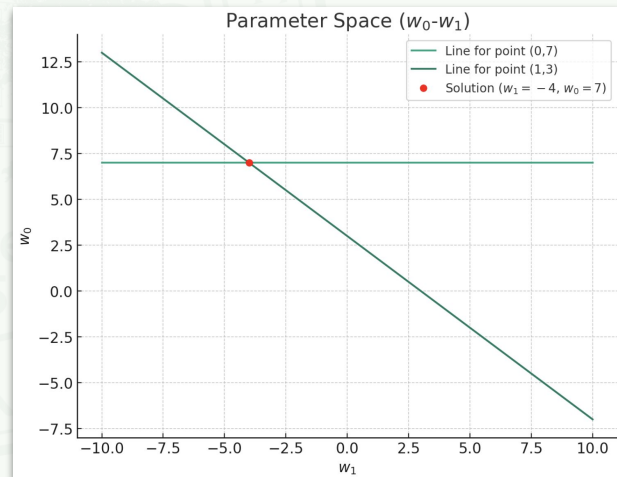Let's find the line using linear algebra:



$$y = w_0 + w_1 x$$

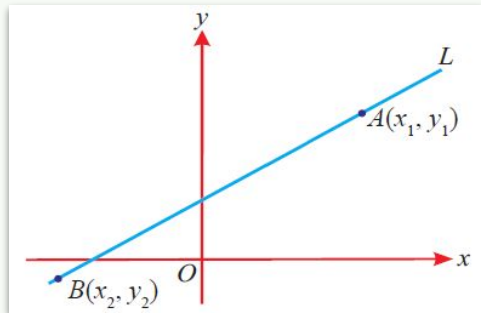Suppose we have these two points:

$$(0, 7), (1, 3)$$

We can substitute these points into two of our model equations to get two **relations between the weights**.

# We don't want to search for intersections of lines!

Let's find the line using linear algebra:



$$y = w_0 + w_1 x$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{bmatrix} 7 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 7 \\ 3 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$
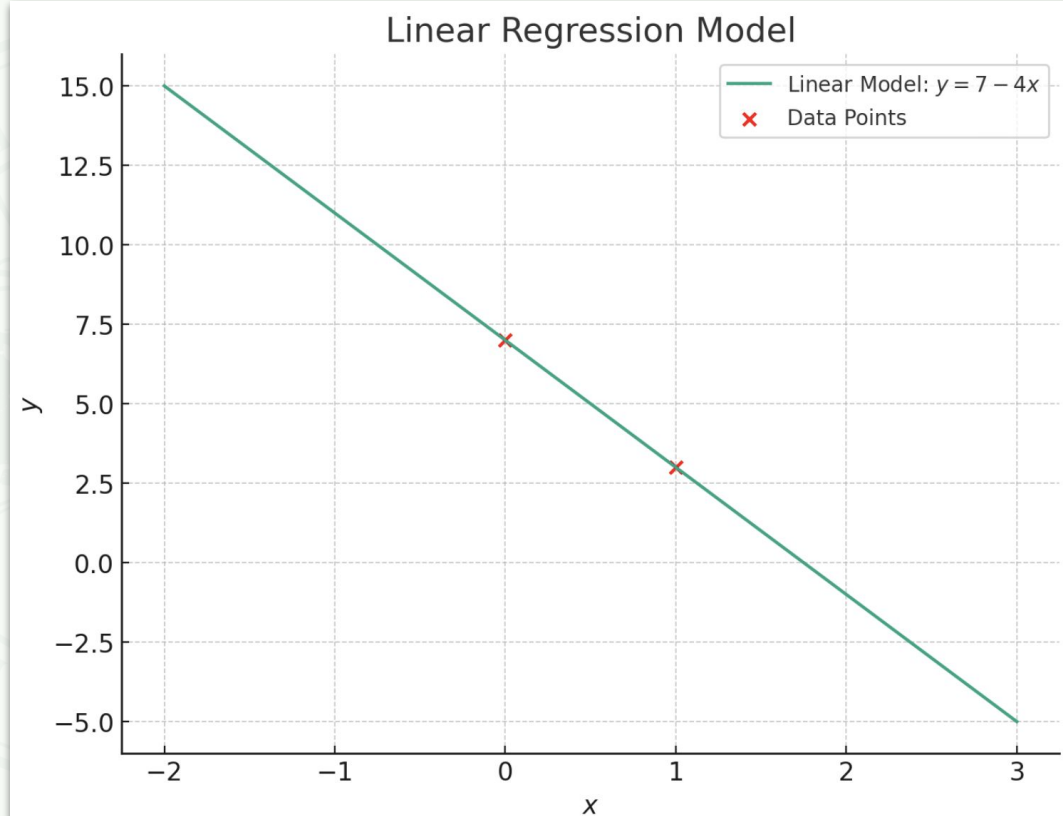
$$\begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 7 \\ 3 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} 7 \\ -4 \end{bmatrix}$$

$$y = 7 - 4x$$

CMSE

# We have our final model!

$$\begin{bmatrix} 7 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 7 \\ 3 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 7 \\ 3 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} 7 \\ -4 \end{bmatrix}$$

$$y = 7 - 4x$$



Linear Regression Model

Legend: Linear Model: $y = 7 - 4x$; Data Points (×)

CMSE

# Key Step: Finding the Inverse

The inverse is defined through:

$$XX^{-1} = I$$

For a 2x2 matrix:

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$\det(A) = ad - bc$$

By hand:

$$\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} X^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Note that the determinant must not be zero!

$$a = 1$$

$$b = 0$$

$$a + c = 0$$

$$b + d = 1$$

CMSE

1. there is no solution

2. there are infinitely many solutions

# Understanding Matrix Rank in Linear Algebra

> **Rank of a Matrix**: The rank of a matrix is the dimension of the vector space spanned by its rows (row rank) or columns (column rank). It is equal to the maximum number of linearly independent row vectors or column vectors in the matrix.

- **Rank Deficient**: A matrix is rank deficient if its rank is less than its number of rows or columns. This implies that some rows or columns are linear combinations of others.

    a. **Geometric Interpretation**: Rank deficiency means that the space spanned by the row or column vectors of the matrix is "flattened" into a lower dimension than the matrix itself, indicating dependency among the vectors.

- **Inconsistent System**: A system of linear equations is inconsistent if there are no solutions that satisfy all equations simultaneously.

    a. **Geometric Interpretation**: Geometrically, an inconsistent system represents a set of planes (or lines in two dimensions) that do not intersect at a common point. Instead, at least two of them are parallel, which means there is no point that lies on all planes (or lines).

$$\begin{bmatrix} 1 & 8 & 13 & 12 \\ 14 & 11 & 2 & 7 \\ 4 & 5 & 16 & 9 \\ 15 & 10 & 3 & 6 \end{bmatrix}$$

CMSE

# Python Libraries

CMSE

# Linear Algebra For Image Compression



Original data

CMSE

# Bias-Variance Tradeoff

How should we pick our model?

- a line/hyperplane?

- radial basis function?

  - which one?

> "Model selection" is the process of selecting the best model among a set of models using data.



High variance — overfitting

High bias — underfitting

Low bias, low variance — Good balance

# Real Data: Too Many Equations, Not Enough Unknowns!



$6x + 2y = 10$

$(-2, 11)$
$(-1, 8)$
$(0, 5)$
$(1, 2)$
$(2, -1)$

unrealistic!

realistic

CMSE

# Data Matrix

- columns are features - perhaps as many as hundreds

- rows are data measurements - could have tens of thousands, or more

| Observation Number | Temperature ($x_i$) | Yield ($y_i$) |
|---|---|---|
| 1 | 50 | 122 |
| 2 | 53 | 118 |
| 3 | 54 | 128 |
| 4 | 55 | 121 |
| 5 | 56 | 125 |
| 6 | 59 | 136 |
| 7 | 62 | 144 |
| 8 | 65 | 142 |
| 9 | 67 | 149 |
| 10 | 71 | 161 |
| 11 | 72 | 167 |
| 12 | 74 | 168 |
| 13 | 75 | 162 |
| 14 | 76 | 171 |
| 15 | 79 | 175 |
| 16 | 80 | 182 |
| 17 | 82 | 180 |
| 18 | 85 | 183 |
| 19 | 87 | 188 |
| 20 | 90 | 200 |
| 21 | 93 | 194 |
| 22 | 94 | 206 |
| 23 | 95 | 207 |
| 24 | 97 | 210 |
| 25 | 100 | 219 |

CMSE

# Linear Algebra for Realistic Data Matrices



$$y = w_0 + w_1 x_1 + w_2 x_2,$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & & \\ 1 & x_{1N} & x_{2N} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

$\underbrace{\phantom{xxx}}_{N \times 1}$  $\underbrace{\phantom{xxxxxx}}_{N \times 3}$  $\underbrace{\phantom{xx}}_{3 \times 1}$

We cannot invert this to find the weights because the data matrix $X$ is not square.

$$\mathbf{y} = X\mathbf{w}$$

CMSE

# We Solved This Previously Using Optimization

$$L(w_0, w_1, w_2) = \sum_d \left(y_d - (w_0 + w_1 x_{1d} + w_2 x_{2d})\right)^2,$$

$$\frac{\partial L}{\partial w_0} = 0,$$

$$\frac{\partial L}{\partial w_1} = 0,$$

$$\frac{\partial L}{\partial w_2} = 0$$

Minimize the distance from the hyperplane to the data points.

It may feel like we have two, very-different methods.

Can we solve the realistic case with linear algebra?

CMSE

# Matrix Times Its Transpose

$$\underbrace{X}_{n \times m} \underbrace{X^T}_{m \times n} = \underbrace{A}_{n \times n},$$

$$\underbrace{X^T}_{m \times n} \underbrace{X}_{n \times m} = \underbrace{B}_{m \times m}$$

Because matrices don't commute, $A$ and $B$ are not the the same matrix. Order matters.

Note that:

1. $A$ and $B$ are square, and therefore might have an inverse.

2. $A$ and $B$ are symmetric matrices.

CMSE

$$\mathbf{y} = X\mathbf{w},$$

$$X^T\mathbf{y} = X^T X\mathbf{w},$$

$$X^T\mathbf{y} = (X^T X)\mathbf{w},$$

$$(X^T X)^{-1}X^T\mathbf{y} = (X^T X)^{-1}(X^T X)\mathbf{w},$$

$$(X^T X)^{-1}X^T\mathbf{y} = I\mathbf{w},$$

$$\boxed{\mathbf{w} = (X^T X)^{-1}X^T\mathbf{y}}$$

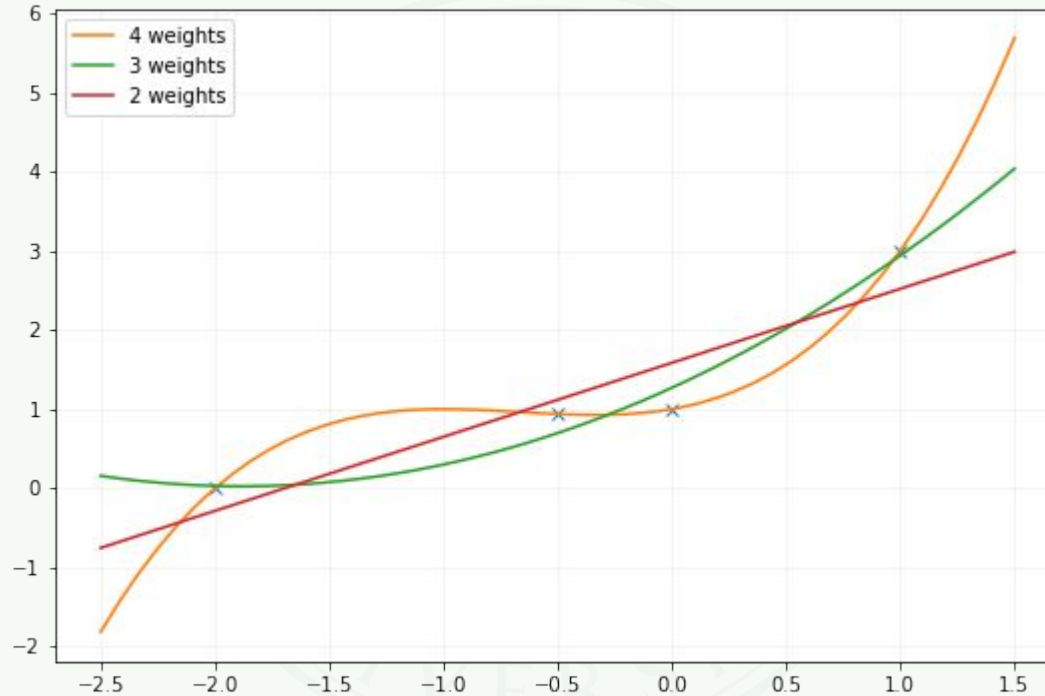**Most important equation in data science?**

CMSE

# Pseudo-Inverse

A generalization of the inverse to non-square matrices is:

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

In this form, $X^+$ is referred to as the "*Moore-Penrose inverse*".

CMSE

# Homework Solution: Polynomial Regression

CMSE

# Matrix Decompositions: Big Picture

It is often convenient to write a matrix as a product of other matrices.

This might seem like a step backwards, but there are very good reasons for doing this.

Random example:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} \ell_{11} & 0 & 0 \\ \ell_{21} & \ell_{22} & 0 \\ \ell_{31} & \ell_{32} & \ell_{33} \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix}$$

$$A = LU$$

This decomposition is often used to find the inverse and/or determinant of a matrix. You use this when you use `linalg`.

CMSE

# Data Science: Singular Value Decomposition (SVD)

There is one decomposition that occurs across most of data science: the SVD. We'll focus on the SVD for this reason.

The SVD is important for two reasons:

- it generalizes ideas for square matrices to non-square data matrices,

- it reveals structure in the data that can be exploited.

$$A = U\Sigma V^T$$

CMSE

# SVD: Details

$$A = U \Sigma V^T$$

Dimensions: $A$ is $m \times n$, $U$ is $m \times m$, $\Sigma$ is $m \times n$, $V^T$ is $n \times n$.

$U$ and $V$ are orthogonal (more generally: unitary):

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

$$U^T U = I$$

$\Sigma$ is a rectangular diagonal matrix:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

CMSE

# Singular Values

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

The diagonal values of Σ are called "singular values", and by convention are taken to be ordered from largest to smallest:

$$\sigma_1 > \sigma_2 > \sigma_3 > \ldots$$

Recall from last week that the columns of a matrix form a basis in a subspace. The number of linearly independent columns is called the "rank" of the matrix.
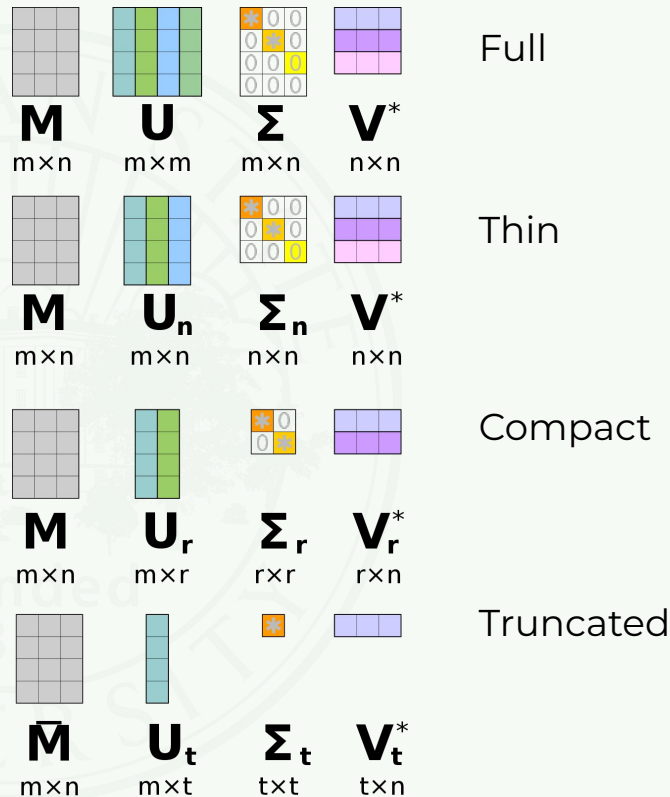
The number of non-zero singular values {$\sigma_n$} is the rank of *A*.

CMSE

# Reduced SVDs: Thin, Compact, Skinny....

Sometimes you will see the SVD in a reduced form.

Usually this is for computational reasons. For example, reduced SVDs can save memory.

Often you can tell the library which form you want returned.



Full

$$\mathbf{M} \qquad \mathbf{U} \qquad \mathbf{\Sigma} \qquad \mathbf{V}^*$$
$$m{\times}n \qquad m{\times}m \qquad m{\times}n \qquad n{\times}n$$

Thin

$$\mathbf{M} \qquad \mathbf{U_n} \qquad \mathbf{\Sigma_n} \qquad \mathbf{V}^*$$
$$m{\times}n \qquad m{\times}n \qquad n{\times}n \qquad n{\times}n$$

Compact

$$\mathbf{M} \qquad \mathbf{U_r} \qquad \mathbf{\Sigma_r} \qquad \mathbf{V}_r^*$$
$$m{\times}n \qquad m{\times}r \qquad r{\times}r \qquad r{\times}n$$

Truncated

$$\mathbf{\bar{M}} \qquad \mathbf{U_t} \qquad \mathbf{\Sigma_t} \qquad \mathbf{V}_t^*$$
$$m{\times}n \qquad m{\times}t \qquad t{\times}t \qquad t{\times}n$$

CMSE

# Linear Regression and Pseudoinverse Via SVD

$$\mathbf{y} = K\mathbf{w},$$

(now, $K$ is not square - we can't easily invert this)

$$= U\Sigma V^T\mathbf{w},$$

need to define this

$$U^T\mathbf{y} = \Sigma V^T\mathbf{w},$$

$$\Sigma^{-1}U^T\mathbf{y} = V^T\mathbf{w},$$
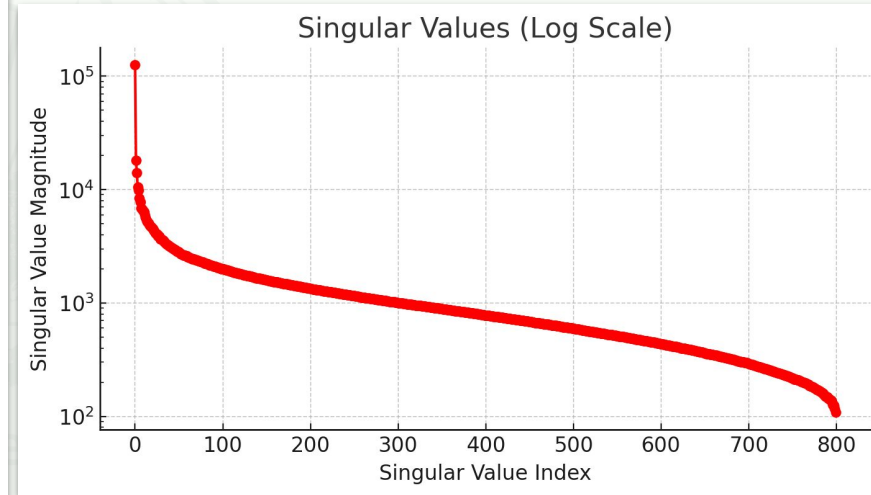
$$\mathbf{w} = V\Sigma^{-1}U^T\mathbf{y},$$

$$K^+ = V\Sigma^+ U^T$$

$\Sigma^+$ is the pseudoinverse of Σ, formed by replacing all non-zero entries by their reciprocal and transposing.

All zero entries remain zero.

CMSE

# Let's Play With SVD



Original Grayscale Image



Singular Values (Log Scale)

CMSE

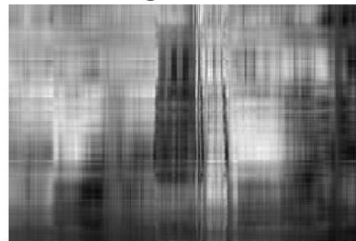# Let's Play With SVD



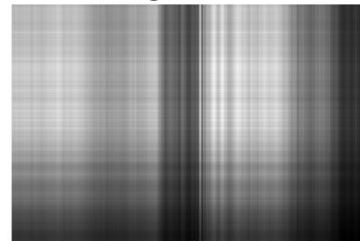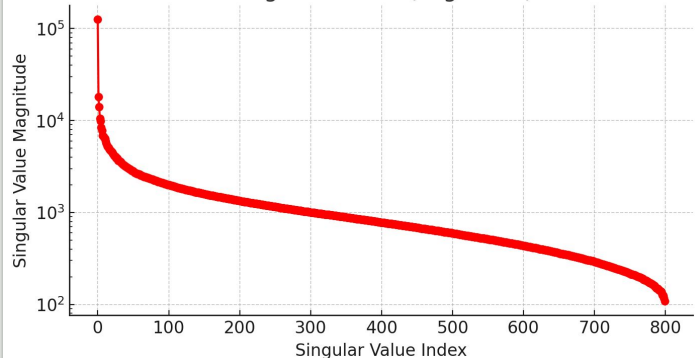800 Singular Values | 50 Singular Values | 10 Singular Values | 5 Singular Values | 1 Singular Values
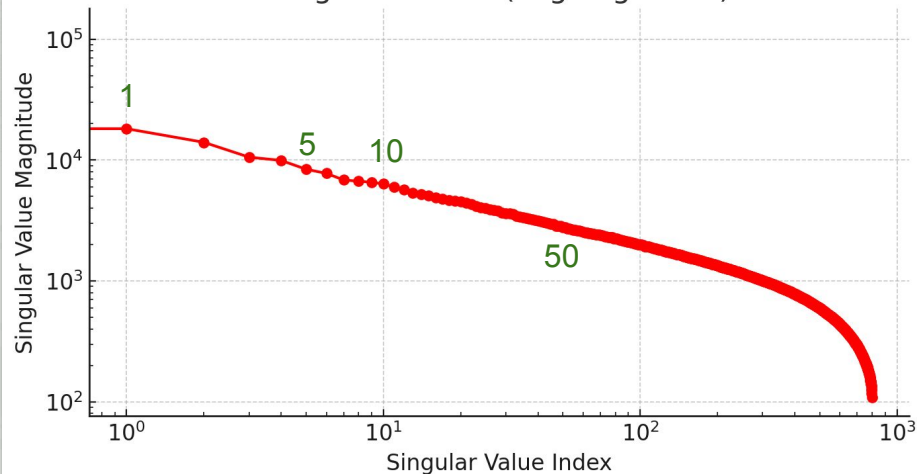


Singular Values (Log Scale)

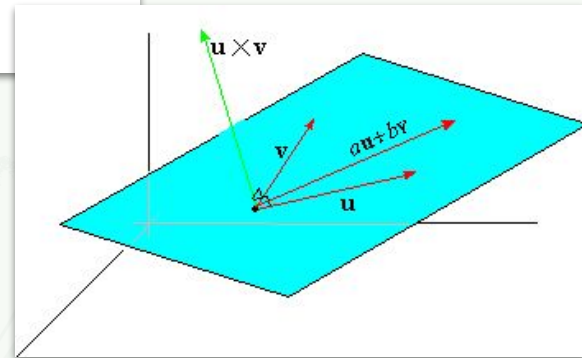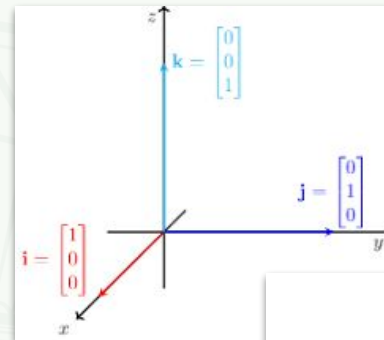Singular Values (Log-Log Scale)

CMSE

# See you Wednesday!

CMSE

# Column Space: Geometric Interpretation

$$\mathbf{y} = K\mathbf{w},$$

$$= \begin{bmatrix} | & | & \dots & | \\ K_1 & K_2 & \dots & K_n \\ | & | & \dots & | \end{bmatrix} \mathbf{w},$$

$$= w_0 K_1 + w_1 K_2 + \dots + w_{n-1} K_n$$





Each column acts as a "unit vector". If **y** is not in the "span" of the columns, then there is no solution.

CMSE