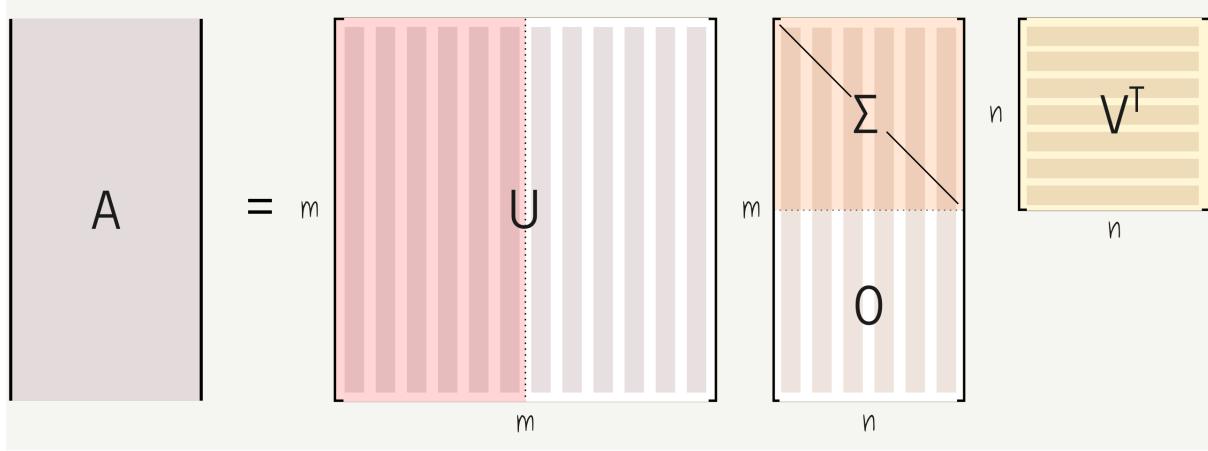


# Understanding Linear Regression using the Singular Value Decomposition

[ machine-learning (/tag/machine-learning) deep-learning (/tag/deep-learning)  
 linear-regression (/tag/linear-regression) svd (/tag/svd)  
 singular-value-decomposition (/tag/singular-value-decomposition) ]



Oct 12, 2020



([https://twitter.com/intent/tweet?  
 text=Understanding Linear Regression using the  
 Singular Value  
 Decomposition&url=https://sthalles.github.io/svd-for-regression/  
 &via=jekyllrb&related=jekyllrb](https://twitter.com/intent/tweet?text=Understanding%20Linear%20Regression%20using%20the%20Singular%20Value%20Decomposition&url=https://sthalles.github.io/svd-for-regression/&via=jekyllrb&related=jekyllrb))



([https://facebook.com/sharer.php?  
 u=https://sthalles.github.io/svd-for-regression/](https://facebook.com/sharer.php?u=https://sthalles.github.io/svd-for-regression/))



([https://plus.google.com/share?  
 url=https://sthalles.github.io/svd-for-regression/](https://plus.google.com/share?url=https://sthalles.github.io/svd-for-regression/))



([http://www.linkedin.com/shareArticle?  
 mini=true&url=https://sthalles.github.io/svd-for-regression/&title=Understanding Linear  
 Regression using the Singular Value  
 Decomposition&summary=&source=<DOMAIN>](http://www.linkedin.com/shareArticle?mini=true&url=https://sthalles.github.io/svd-for-regression/&title=Understanding%20Linear%20Regression%20using%20the%20Singular%20Value%20Decomposition&summary=&source=<DOMAIN>))

## Introduction

It is very common to see blog posts and educational material explaining linear regression. In most cases, probably because of the big data and deep learning biases, most of these educational resources take the gradient descent approach to fit lines,

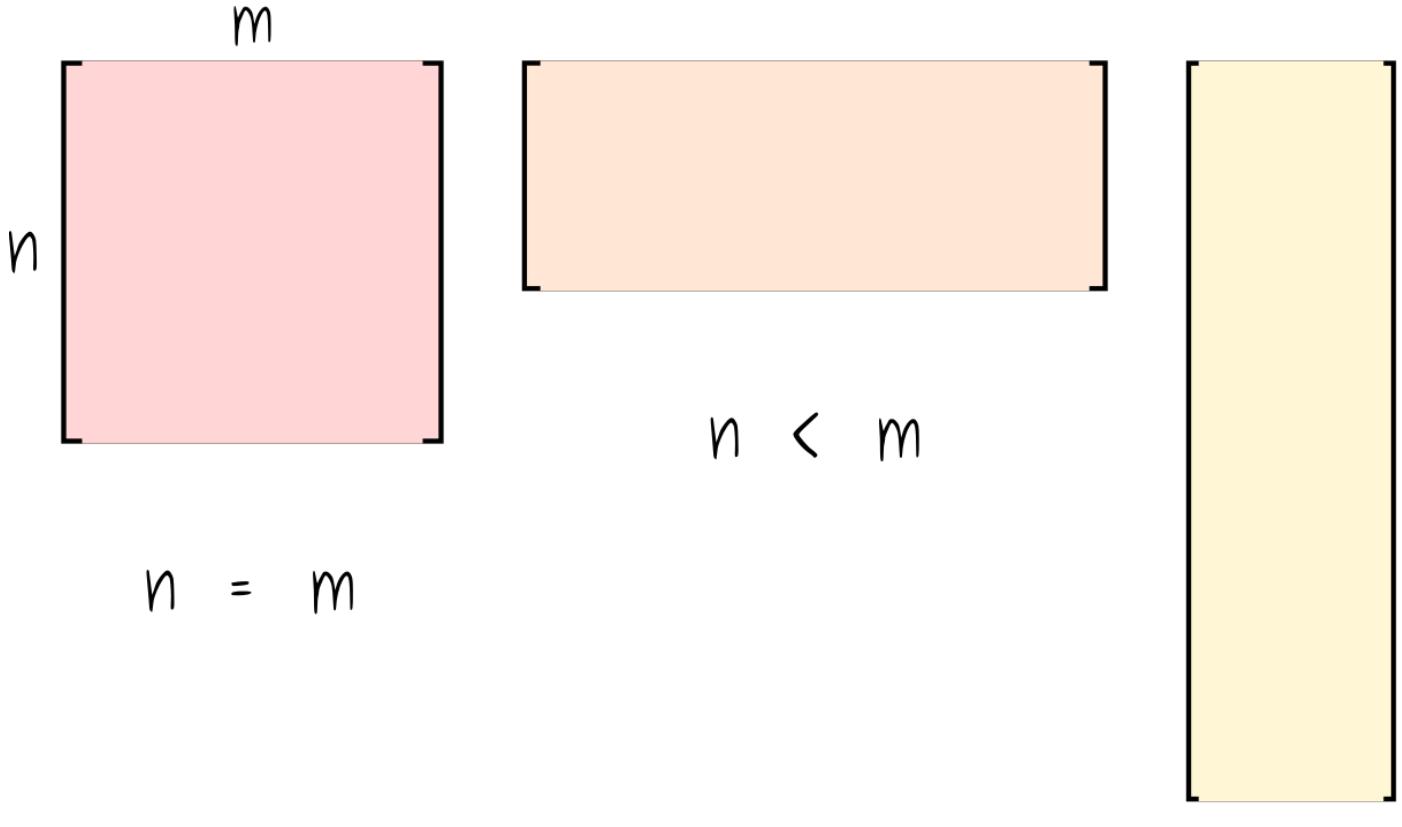
planes, or hyperplanes to high dimensional data. In this post, we will also talk about solving linear regression problems but through a different perspective. Most specifically, we will talk about one of the most fundamental applications of linear algebra and how we can use it to solve regression problems. Yes, I am talking about the SVD or the Singular Value Decomposition. This computational tool is used as a basis to solve a myriad of problems, including dimensionality reduction, with PCA, and statistical learning using linear regression.

## Linear Models and Systems of Linear Equations

Through the lens of linear algebra, a regression problem reduces to solving systems of linear equations of the form  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Here,  $\mathbf{A}$  and  $\mathbf{b}$  are known, and  $\mathbf{x}$  is the unknown. We can think of  $\mathbf{x}$  as our model. In other words, we want to solve the system for  $\mathbf{x}$ , and hence,  $\mathbf{x}$  is the variable that relates the observations in  $\mathbf{A}$  to the measures in  $\mathbf{b}$ .

Here,  $\mathbf{A}$  is a data matrix. We can think of the rows of  $\mathbf{A}$  as representing different instances of the same phenomenon. They can represent records for individual patients submitted to a hospital, records for different houses being sold, or pictures of different people's faces. Complementary, we can view the columns of the matrix  $\mathbf{A}$  as recording different characteristics of each instance in the rows of  $\mathbf{A}$ . In a patient hospital example, such features might include the blood pressure when he/she arrived at the hospital or if the patient has had a surgical procedure or not.

Also, note that the matrix  $\mathbf{A}$  might have different shapes. First,  $\mathbf{A}$  could be a square matrix. Yes, it is very unlikely (for the situations we usually encounter in data science) but otherwise possible.



$n > m$

The matrix  $A$  may have different shapes. It can be squared. It can be wider and short, or it can be tall and skinny.

Second,  $A$  could have more columns than rows. In this scenario,  $A$  would have a short and wide shape. And lastly, (and that is the most usual case in data science), the matrix  $A$  assumes the form of a tall and skinny matrix, with many more rows than columns.

But why should I care for the shape of the matrix  $A$ ?

Interestingly, the shape of  $A$  will dictate whether the linear system of equations has a solution, has infinitely many solutions, or does not have a solution at all.

Let's start with the boring case. If the matrix is squared (number of rows equals the number of columns) and it is invertible, meaning that the matrix  $A$  has full rank (all columns are linearly independent), that pretty solves the problem.

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

If the matrix  $A$  is squared and invertible, the system of equations has a solution.

However, if the matrix has more columns than it has rows, we are likely dealing with the case where there are infinitely many solutions. To visualize this curious scenario, picture a  $3 \times 6$  matrix, i.e., 3 rows and 6 columns. We can think of it as having a 3D space and 6 different vectors that we can use to span the 3D space. However, to span a 3D space, we only need 3 linearly independent vectors, but we have 6! This leaves 3 dependent vectors that can be used to formulate infinitely many solutions.

Finally, by analogy, if we have a matrix  $A$  with more rows than columns, we can view it as trying to span a very high-dimensional space with fewer vectors than we would need. For instance, picture a matrix with 6 rows and 2 columns. Here, we have a 6D space, but we only got 2 vectors to span it. It does not matter how much we try it, in the best case, we can only span a plane on 6D. And that is crucial because we only have a solution to  $Ax = b$  if the vector  $b$  is in the column space of  $A$ . But here, the column space of  $A$  spans 2D subspace (a plane) on a much larger 6D space. This makes the probability of the vector  $b$  to be in the subspace spanned by the columns of  $A$  improbable.

To visualize how unlikely it is, picture a 3D space and a subspace spanned by two vectors (a plane in 3D). Now, imagine you choose 3 values at random. This will give you a point on the 3D space. Now, ask yourself: what is the probability that my randomly chosen point will be on the plane?

Nonetheless, in situations where we do not have a solution for a linear system of equations  $Ax = b$  (or we have infinitely many solutions), we still want to do our best. And to do this, we need to find the best approximate solution. Here is where the SVD kicks in.

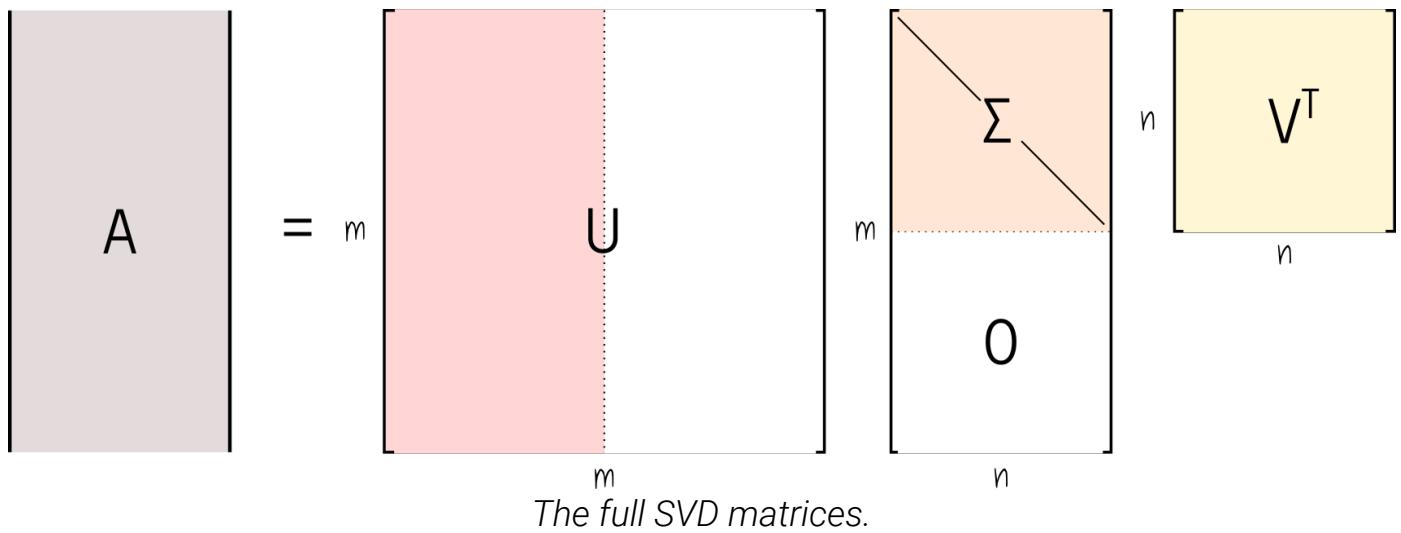
## A Short Intro to the SVD

The main idea of the singular value decomposition, or SVD, is that we can decompose a matrix  $A$ , of any shape, into the product of 3 other matrices.

$$A = USV^T$$

Given a matrix of any shape, the SVD decomposes  $A$  into a product of 3 matrices:  $U$ ,  $\Sigma$ ,  $V^T$ .

Here,  $U$  is an  $m \times m$  square matrix,  $\Sigma$  is a rectangular matrix of shape  $m \times n$ , and  $V^T$  is a square matrix and has shape  $n \times n$ .



The matrices  $U$  and  $V^T$  have a very special property. They are **unitary matrices**. One of the main benefits of having unitary matrices like  $U$  and  $V^T$  is that if we multiply one of these matrices by its transpose (or the other way around), the result equals the identity matrix.

On the other hand, the matrix  $\Sigma$  is diagonal, and it stores non-negative singular values ordered by relevance.

$$\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$$

$$\mathbf{V}\mathbf{V}^T = \mathbf{V}^T\mathbf{V} = \mathbf{I}$$

*Property of unitary matrices.*

Note that, since the  $\Sigma$  matrix is diagonal, only the first  $n$  row diagonal values are worth keeping. Indeed the last  $n$  rows of  $\Sigma$  are filled with 0s. For this reason, it is very common to keep only the first  $r \times r$  non-negative diagonal values of  $\Sigma$ , along with the corresponding  $r$  columns and rows of  $\mathbf{U}$  and  $\mathbf{V}^T$  respectively. Note that  $r = \min(m, n)$ . This is commonly referred to as the economy (or compact) SVD, and from this point on, we will assume the matrices  $\mathbf{U}$ ,  $\Sigma$ , and  $\mathbf{V}^T$  are derived from the economy procedure.

Quick note, it is very common to also truncate the SVD based on some criteria. Under some assumptions, it is possible to find an optimal threshold for suppressing some of the singular values with small magnitudes.

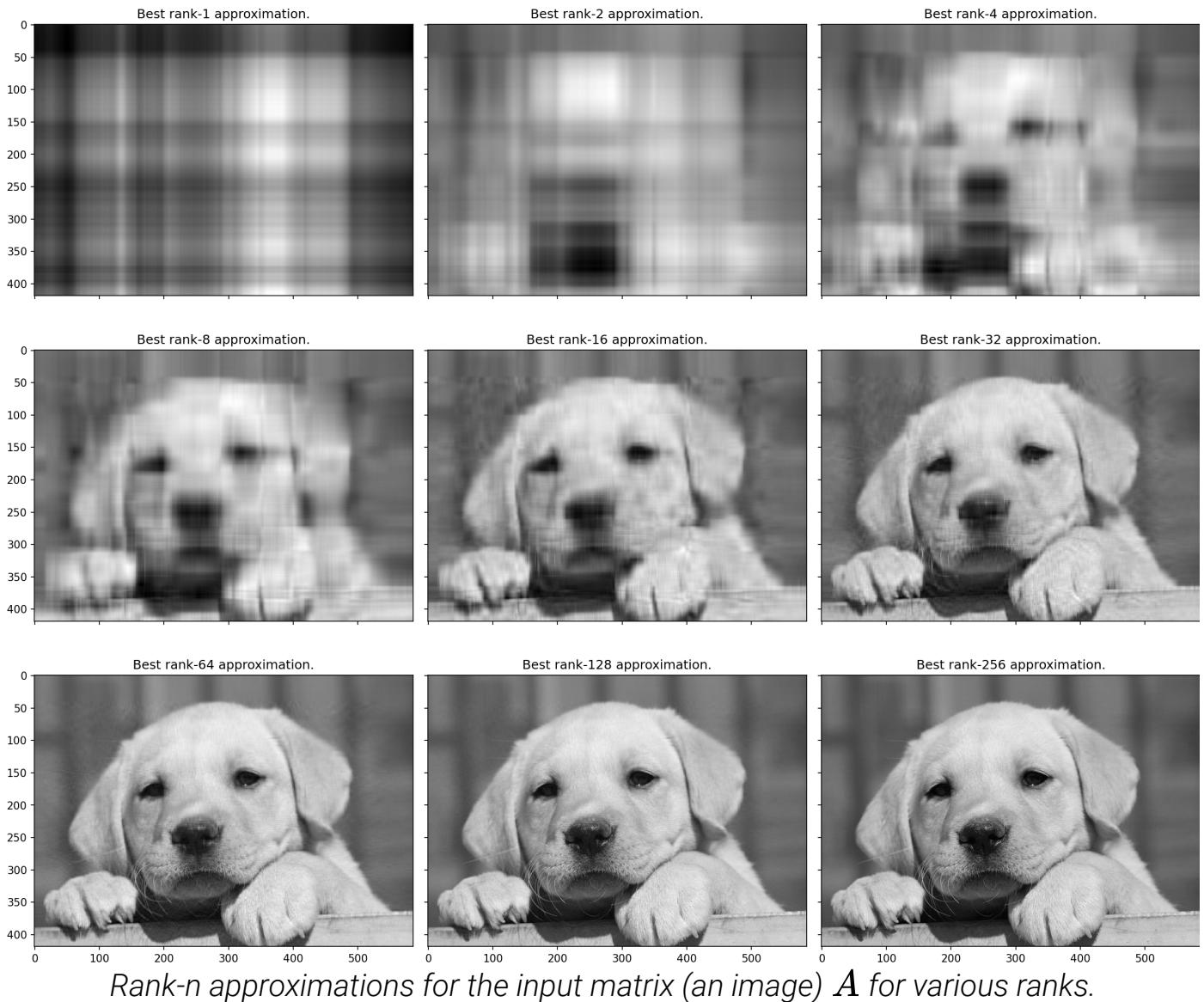
$$\mathbf{A} = m \begin{bmatrix} \mathbf{U} \end{bmatrix} \begin{bmatrix} \Sigma \end{bmatrix} \begin{bmatrix} \mathbf{V}^T \end{bmatrix}$$

*The economy SVD data matrices.*

It is important to note that the economy SVD produces a change in the shape of the matrices  $\mathbf{U}$  and  $\Sigma$  (if one of the diagonal values of  $\Sigma$  is zero,  $\mathbf{V}^T$  also suffers a shape change). If the diagonal values of  $\Sigma$  are all positives, thus  $r = n$ , we discard the right

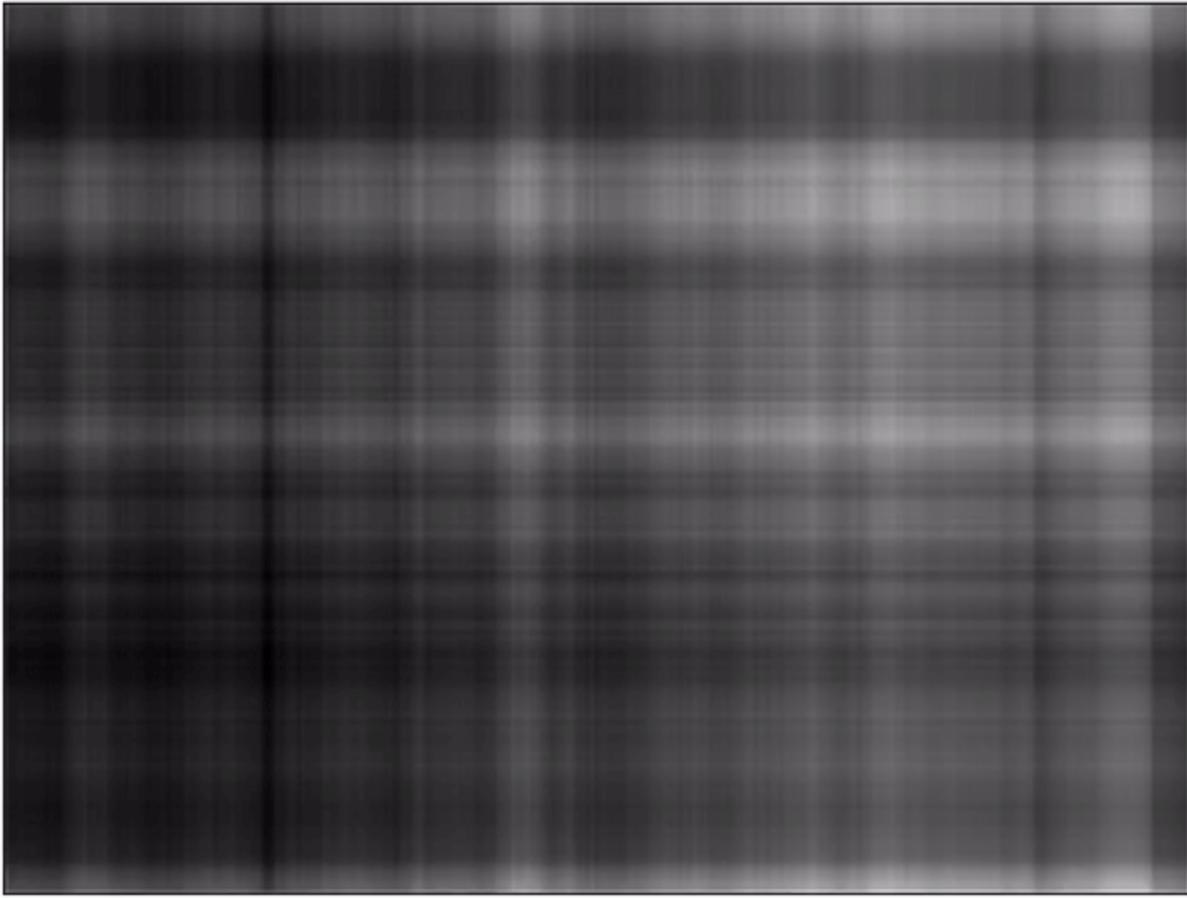
half of the  $\mathbf{U}$  matrix (the orthogonal complement of  $\mathbf{U}$ ), which gives  $\mathbf{U}$  a rectangular  $m \times r$  shape. More critical,  $\mathbf{U}$ , and possibly  $\mathbf{V}^T$ , are now semi-unitary matrices, which means that only  $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}$ .

The SVD provides a basis that allows us to reconstruct the input signal in terms of low-rank matrix approximations. Let me be more clear. If we combine each column of  $\mathbf{U}$  with the corresponding row of  $\mathbf{V}^T$ , and scale the resulting matrix by the corresponding  $\sigma$  value, we will get the best rank-1 approximation of  $\mathbf{A}$  in terms of least squares.



*Rank-n approximations for the input matrix (an image)  $\mathbf{A}$  for various ranks.*

And as we continue combining the columns of  $\mathbf{U}$  with rows of  $\mathbf{V}^T$ , scaled by the corresponding  $\sigma$ , we get the next best rank-i approximation of the data matrix  $\mathbf{A}$ . Indeed, that is another excellent application of the SVD—data compression. But that is a subject for another writing.

**Best rank-1 approximation.**

*Image reconstruction using SVD's  $\mathbf{U}$ ,  $\Sigma$  and  $\mathbf{V}^T$  matrices. Using the first 256 singular values, we get the best rank-256 approximation of the original input image.*

As we said before, the problem of working with a non-square matrix  $\mathbf{A}$  is that we cannot invert it. And that is the main reason why we cannot solve the system of equations as we would for a square matrix  $\mathbf{A}$ . However, if we cannot invert the matrix  $\mathbf{A}$ , I invite you to ask yourself the following question.

**What would be the best matrix  $\mathbf{A}^+$  that, when multiplied by  $\mathbf{A}$ , would come as close as possible to identity matrix  $\mathbf{I}$ ?**

The answer to this question solves the problem of finding the best possible solution when the system of equations has infinitely many solutions or no solution at all. Luckily, the answer also lies in the SVD.

If we know that the SVD always exists (for matrices of any shape), and by combining the columns of  $\mathbf{U}$ , the rows of  $\mathbf{V}^T$ , and the singular values  $\sigma$ , we can reconstruct the original input matrix nearly perfectly, what happens if we try to invert the SVD?

Let me spoil it with no further ado. It turns out that the best matrix  $\mathbf{A}^+$  that approximately solves the question  $\mathbf{A}^+ \mathbf{A} \approx \mathbf{I}$  is the inverse of the SVD. In other words, the best approximation for  $\mathbf{A}^{-1}$  is  $\mathbf{S}\mathbf{V}\mathbf{D}^{-1}$ . Let's follow the math.

$$Ax = b$$

$$USV^T x = b$$

$$(USV^T)^{-1}(USV^T)x = (USV^T)^{-1}b$$

$$VS^{-1}U^T USV^T x = VS^{-1}U^T b$$

$$\hat{x} = VS^{-1}U^T b$$

$$\hat{x} = A^+ b$$

$$\# A^+ = VS^{-1}U^T$$

Finding the pseudo-inverse of  $A$  through the SVD. The pseudo-inverse  $A^+$  is the closest we can get to non-existent  $A^{-1}$

First, we compute the SVD of  $A$  and get the matrices  $USV^T$ . To solve the system of equations for  $x$ , I need to multiply both sides of the equation by the inverse of the SVD matrices. Luckily now, it is very easy to invert each one of the 3 SVD matrices. To invert the product of the 3 matrices  $USV^T$ , I take the product of the inverse matrices in reverse order!

After inverting the matrices, if we look closely at the left-hand side, we can see that most matrices will cancel like crazy, leaving us with the best approximate solution for  $\hat{x}$ . Note that since the matrix  $U$  is semi-unitary, only  $U^T U = I$  holds. Moreover, if (and we assume that) all the singular values are non-negative, then  $V^T$  continuous to be a unitary matrix. Hence, to invert  $U$  and  $V^T$ , we just multiply each one by their transpose, i.e.,  $U^T U = I$  and  $V V^T = I$ .

$$Ax$$

$$USV^T VS^{-1} U^T b$$

$$UU^T b$$

Finding the projection of  $b$

If we go further and substitute our best solution  $\hat{x}$  into  $A\hat{x}$ , we will see that most of the matrices cancel each other as well, until we reach  $UU^T$ . As we said before,  $U$  is semi-unitary, and  $UU^T$  is not the identity matrix. Instead,  $UU^T$  is the projection of  $b$  onto the

subspace spanned by the columns of  $\mathbf{U}$  (hence columns of  $\mathbf{A}$ ), which is the best approximate solution in terms of least squares, i.e., we found the least squares solution  $\hat{\mathbf{x}} = \text{minimum}(|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}|_2)$ .

Note that if  $\mathbf{A}$  has more columns than rows and infinitely many solutions, the SVD picks the solution with the minimum 2-norm, i.e.,  $\hat{\mathbf{x}} = \text{minimum}(|\hat{\mathbf{x}}|_2)$ .

## Linear Regression with the SVD

Once we have established the required SVD jargon, we can use it to find approximate solutions for real-world problems. In this example, I am going to use the [Boston house-prices dataset](#) ([https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\\_boston.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_boston.html)). The house-prices data matrix  $\mathbf{A}$  contains 506 rows (representing individual houses), and 13 columns (each describing a different characteristic of the houses). Some of these 13 features include:

- Per capita crime rate by town
- The average number of rooms per dwelling
- Weighted distances to five Boston employment centers

You can see the [full description here](#) (<https://scikit-learn.org/stable/datasets/index.html#boston-dataset>).

We want to predict the **median value home price in \$1000's**. These measurements are real values ranging from 5 to 50, and they represent the  $\mathbf{b}$  vector in our system of equations  $\mathbf{Ax} = \mathbf{b}$ .

As usual, the matrix has many more rows than columns. This means that we cannot invert  $\mathbf{A}$  to find the solution to  $\mathbf{Ax} = \mathbf{b}$ . Also, it drastically reduces the possibilities of finding a solution. Indeed, such a solution would only be possible if  $\mathbf{b}$  is a linear combination of the columns of  $\mathbf{A}$ . However, using the SVD, we will be able to derive the pseudo-inverse  $\mathbf{A}^+$ , to find the best approximate solution in terms of least squares – **which is the projection of the vector  $\mathbf{b}$  onto the subspace spanned by the columns of  $\mathbf{A}$** .

The code is very simple to follow, and the results are excellent. Indeed, they are the best possible for a linear model.

```

1 # load the boston house data
2 data = load_boston(return_X_y=False)
3
4 # extract the input data matrix and the targets

```

```

5   A = data.data
6   b = data.target
7
8   # append a column of 1s (these are the biases)
9   A = np.column_stack([np.ones(A.shape[0]), A])
10
11  # split the data into train and test sets
12  X_train, X_test, y_train, y_test = train_test_split(A, b, test_size=0.50, random_state=42)
13
14  # calculate the economy SVD for the data matrix A
15  U,S,Vt = np.linalg.svd(X_train, full_matrices=False)
16
17  # solve Ax = b for the best possible approximate solution in terms of least squares
18  x_hat = Vt.T @ np.linalg.inv(np.diag(S)) @ U.T @ y_train
19
20  # perform train and test inference
21  y_pred = X_train @ x_hat
22  test_predictions = X_test @ x_hat
23
24  # compute train and test MSE
25  train_mse = np.mean((train_predictions - y_train)**2)
26  test_mse = np.mean((test_predictions - y_test)**2)
27
28  print("Train Mean Squared Error:", train_mse)
29  print("Test Mean Squared Error:", test_mse)

```

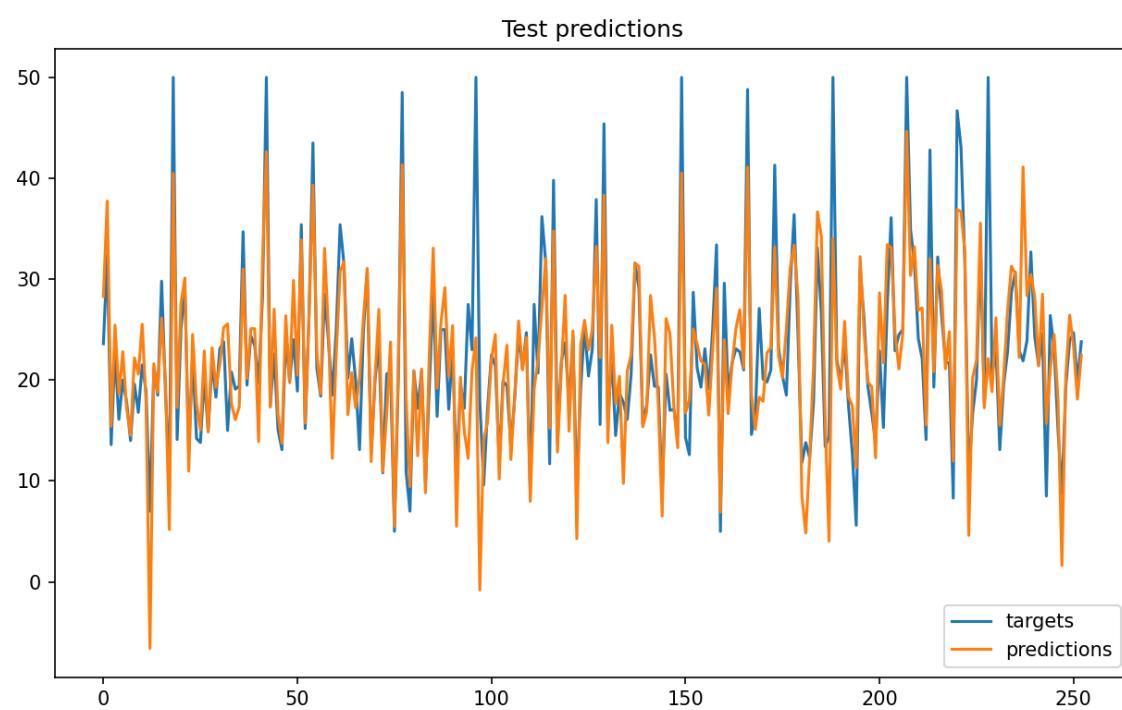
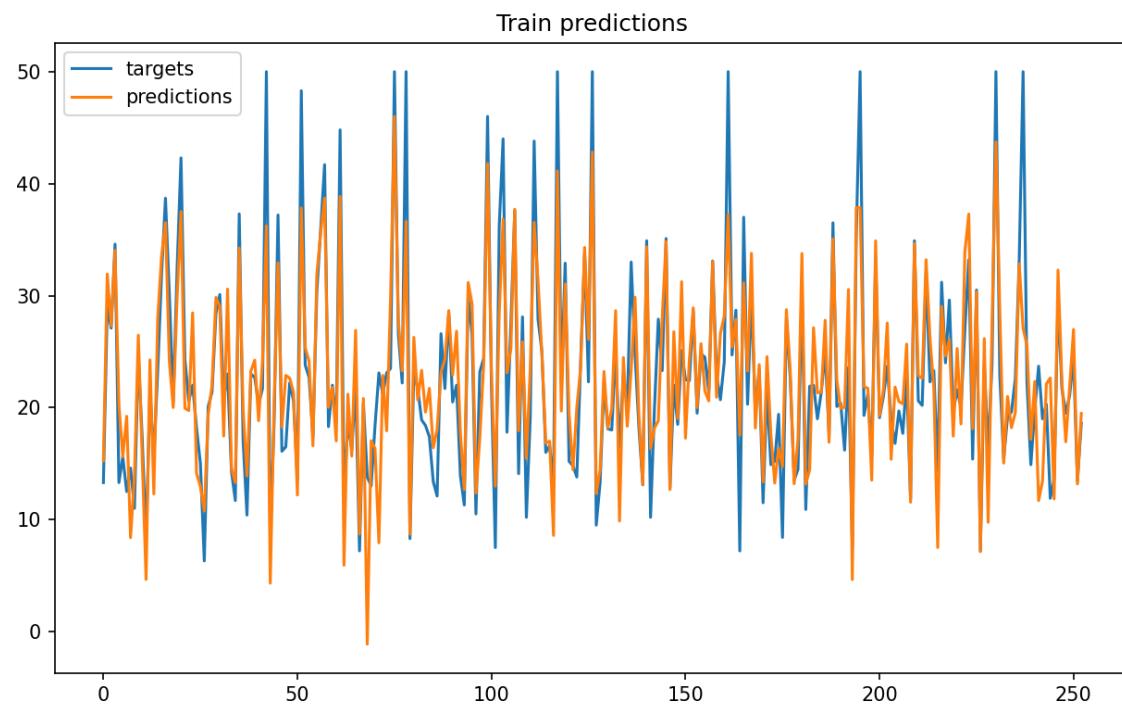
[View raw](#)

[ps://gist.github.com/sthalles/7bbce26fb90da0837758fe4427afabe1/raw/cf9489e5026cda83254d7a54e09b80bed22ceae5/svd-ar-reg.py](https://gist.github.com/sthalles/7bbce26fb90da0837758fe4427afabe1/raw/cf9489e5026cda83254d7a54e09b80bed22ceae5/svd-ar-reg.py)

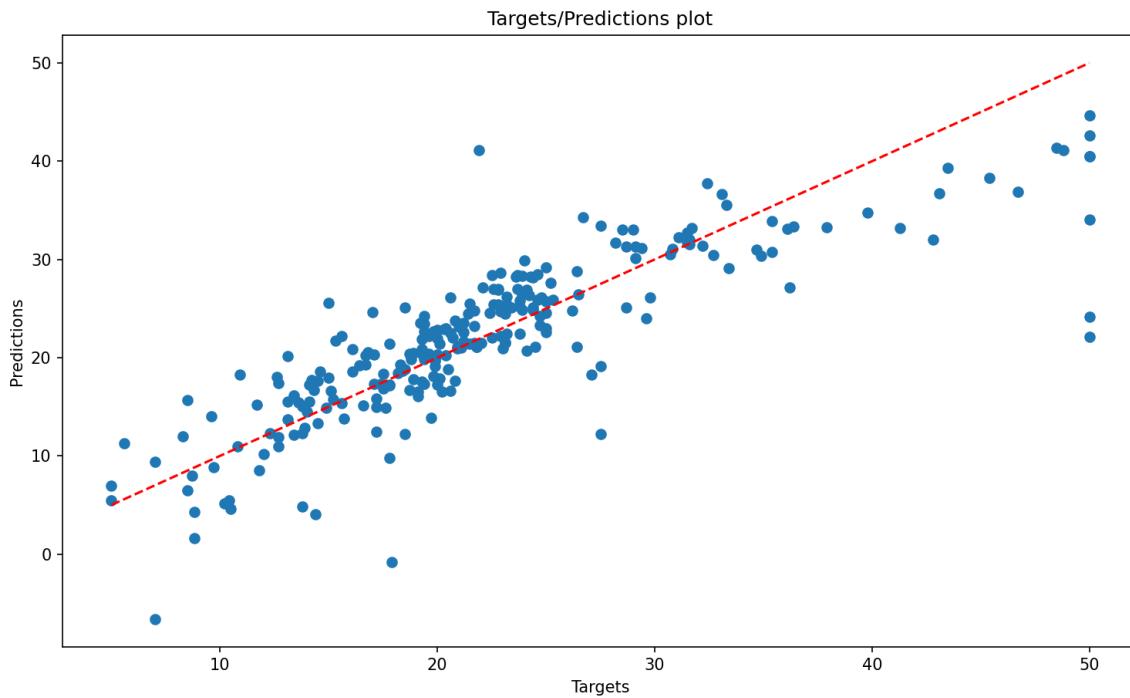
[svd-linear-reg.py](https://gist.github.com/sthalles/7bbce26fb90da0837758fe4427afabe1#file-svd-linear-reg-py) (<https://gist.github.com/sthalles/7bbce26fb90da0837758fe4427afabe1#file-svd-linear-reg-py>) hosted with ❤ by [GitHub](#) (<https://github.com>)

One quick note, look at line 9 of the python code above. At this line, I appended a column full of 1s to the data matrix  $A$ . This column will allow the linear model to learn a bias vector that will add an offset to the hyperplane so that it does not cross the origin.

Take a look at the train and test results below.



*Train and Test predictions for the SVD based linear model.*



*The target/predictions plot allows us to visually assess the correlation between the target values and the model's predictions. Very accurate predictions make the points to be very close to the dotted line.*

Thanks for reading!

Cite as:

```
@article{  
  silva2020svdregression,  
  title={Understanding Linear Regression using the Singular Value Decomposi-  
  author={Silva, Thalles Santos},  
  journal={\url{https://sthalles.github.io}} (https://sthalles.github.io)},  
  year={2020}  
  url={\url{https://sthalles.github.io/svd-for-regression/}} (https://sthalles.gi  
}
```

ALSO ON STHALLES.GITHUB.IO

**How to Add Regularization to ...**

2 years ago • 5 comments

Introduction If you train deep learning models for a living, you might be tired ...

**A Short Introduction to Generative ...**

5 years ago • 1 comment

Let's say there is this very cool party going on at your neighborhood that you ...

**Deeplab Image Semantic ...**

4 years ago • 13 comments

Introduction Deep Convolution Neural Networks (DCNNs) have ...

**Practical Learning**

2 years ago

Introduction Denoising problem

**What do you think?**

12 Responses



Upvote



Funny



Love



Surprised



Angry



Sad

1 Comment

sthalles.github.io

🔒 Disqus' Privacy Policy

1 Login

Heart Favorite 7

Twitter Tweet

f Share

Sort by Best



Join the discussion...

LOG IN WITH

OR SIGN UP WITH DISQUS

Name

Silva, Thalles Santos

thalles753@gmail.com (mailto:thalles753@gmail.com)

Résumé (<https://github.com/sthalles/resume/blob/master/English-Shortened/resume.pdf>)Github (<https://github.com/sthalles>)LinkedIn (<https://www.linkedin.com/in/thalles-silva-32ab08a3/>)