1)

(a) (i) The inverse of a symmetric matrix is itself symmetric.
   **True.**

We want to show $A^{-1} = A^{-T}$ given $A = A^T$

Proof: $I = AA^{-1} = A^T A^{-1}$

multiply $A^{-T}$ on both sides ( left multiply )

$A^{-T} = A^{-T} A^T A^{-1}$

$A^{-T} = A^{-1}$   #

(ii) All $2 \times 2$ orthogonal matrices have the following form

$$\begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} \cos\theta & \sin\theta \\ \sin\theta & -\cos\theta \end{bmatrix} \qquad \underline{\text{True}}$$

let $P \in \mathbb{R}^{2 \times 2}$, $P = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$

if $P$ is orthogonal, then $P^T P = I$

$\Rightarrow \begin{pmatrix} a^2+c^2 & ab+cd \\ ab+cd & b^2+d^2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

$\Rightarrow a^2+c^2 = 1$, $b^2+d^2 = 1$, $ab+cd = 0$

Without loss of generality, ~~we~~ let $a = \cos\theta$, $c = \sin\theta$ and $b = -c$

$$\text{Which are} \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \text{or} \begin{bmatrix} \cos\theta & \sin\theta \\ \sin\theta & -\cos\theta \end{bmatrix} \quad b = \sin\theta \ d = \cos\theta \quad \begin{array}{l} a = d \\ \text{or} \\ b = c, \\ a = -d \end{array}$$

(iii) Assume $A = CC^T$, then we have

$x^T A x = x^T C \cdot C^T x = \|C^T x\|_2^2 \geqslant 0 \Rightarrow A$ is positive semi-definite.

Then $A = U \Sigma U^T$ with $\Sigma_{ii} \geqslant 0$ which means all eigenvalues of $A$ is non-negative.

Also $\operatorname{tr}(A) = \sum_i \lambda_i$   $\lambda_i$ $i=1,2,\cdots$ are the eigenvalues.

$\operatorname{tr}(A) = -8 - 5 - 2 = -15 \Rightarrow$ Not all eigenvalues are non-negative.

Therefore, $A$ cannot be written as $A = CC^T$

2.  (a) (i) $E_Y\left[E_X(X|Y)\right] = E[X]$

proof: $E_Y\left[E_X(Y|X)\right] = \int_{R_y} E[X|y] P(y)\, dy = \int_{R_y}\int_{R_x} x\, P(x|y)\, dx\, P(y)\, dy$

$$= \int_{R_x}\int_{R_y} x\, P(x,y)\, dy\, dx$$

$$= \int_{R_x} x\, P(x)\, dx = E[X]$$

$R_x, R_y$ are the region of $X$ and $Y$ respectively.

(ii) $E\left[I[X\in C]\right] = P(x\in C)$

Proof: $E\left[I[X\in C]\right] = \int_R P(x) I[x\in C]\, dx = \int_C P(x)\, dx = P(x\in C)$

(iii) $var[X] = E_Y\left[var_X[X|Y]\right] + var_Y\left[E_X[X|Y]\right]$

Proof: $var_Y[X|Y] = E_X\left[(X - E[X|Y])^2 | Y\right]$

~~[crossed out text]~~

~~[crossed out text]~~

~~[crossed out text]~~

$E[var(x|Y)] = E\left[E[(x - E[x|Y])^2 | Y]\right]$

$$= E[(x - E[x|Y])^2] \hspace{3cm} \textcircled{1}$$

$var(E_X[x|Y]) = E\left[(E[x|Y] - E(E[x|Y]))^2\right]$

$$= E\left[(E[x|Y] - E[x])^2\right] \hspace{2cm} \textcircled{2}$$

$var(X) = E\left[(x - E[x])^2\right] = E\left[(x - E[x|Y] + E[x|Y] - E[x])^2\right]$

$$= E\left[(x - E[x|Y])^2\right] + E\left[(E[x|Y] - E[x])^2\right]$$
$$+ 2E\left[(x - E[x|Y])\cdot(E[x|Y] - E[x])\right] \hspace{2cm} \textcircled{3}$$

$$E[ (X - E[X|Y]) \cdot ( E[X|Y] - E[X] ) ]$$

$$= E[ X \cdot E[X|Y] - X \cdot E[X] - E[X|Y] \cdot E[X|Y] + E[X] \cdot E[X|Y] ]$$

$$= E[X \cdot E[X|Y]] - \overset{2}{E}[X] - E[ E^2[X|Y] ] + E[X] \cdot E(E[X|Y])$$

$$= E\{ E[X \cdot E[X|Y] | Y] \} - E[ E^2[X|Y] ] - \overset{2}{E}[X] + \overset{2}{E}[X]$$

$$= E\{ E[X|Y] \cdot E[X|Y] \} - E[E^2(X|Y)] = 0 \qquad ④$$

Combining ① ② ③ ④ , we get

$$var(X) = E[ var(X|Y) ] + var( E[X|Y] )$$

(iv) $E[XY] = \displaystyle\iint_{-\infty}^{\infty} xy \, P_{xf}(x,y) \, dx \, dy \qquad$ if $X \perp Y$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \, P_x(x) P_y(y) \, dx \, dy \qquad \text{by independence}$$

$$= \int_{-\infty}^{\infty} x \, P_x(x) \, dx \cdot \int_{-\infty}^{\infty} y \, P_y(y) \, dy$$

$$= E[X] \cdot E[Y]$$

(v) $E[X] = 0 \cdot P(X=0) + 1 \cdot P(X=1) = P(X=1)$

$E[Y] = P(Y=1)$

$E[XY] = \displaystyle\sum_{x=0}^{1} \sum_{y=0}^{1} xy \, P(X=x, Y=y)$

$$= P(X=1, Y=1)$$

$E[XY] = E[X] \cdot E[Y]$

$\Rightarrow P(X=1, Y=1) = P(X=1) \cdot P(Y=1)$
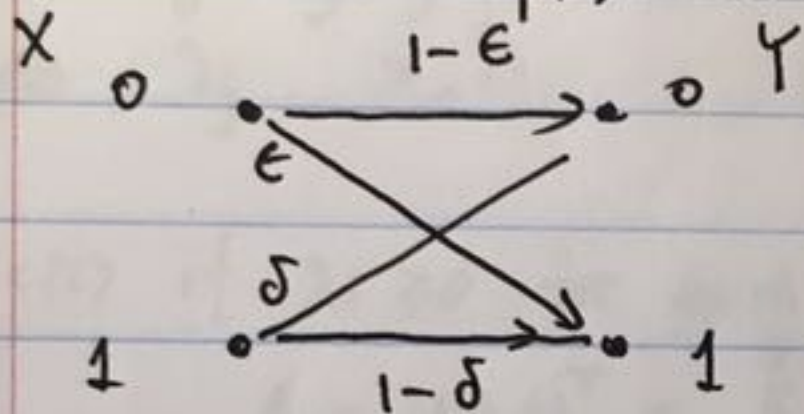
$\Rightarrow X, Y$ are independent.

(b) (i) $P(H=h, D=d) \lneq P(H=h)$

$P(H=h) = \sum_{x} P(H=h, D=x) > P(H=h, D=d)$

(ii) $P(H=h \mid D=d)$ <u>depends</u> $P(H=h)$

$P(H=h) = \sum_{x} P(H=h \mid D=x) P(D=x)$

For example, consider Binary Channel Model in Communication system



$P(Y=0 \mid X=0) = 1-\epsilon$

$P(Y=1 \mid X=0) = \epsilon$

$P(Y=0 \mid X=1) = \delta$

$P(Y=1 \mid X=1) = 1-\delta$

$P(X=0) = \frac{1}{2}$

$P(X=1) = \frac{1}{2}$

Thus $P(Y=0) = P(Y=0 \mid X=0) P(X=0) + P(Y=0 \mid X=1) P(X=1)$

$= (1-\epsilon) \times \frac{1}{2} + \delta \times \frac{1}{2} = \frac{1}{2}(1-\epsilon+\delta)$

$P(Y=0 \mid X=1) = \delta$

The relationship between $P(Y=0)$ and $P(Y=0 \mid X=1)$ depends on the choice of $\epsilon, \delta$.

(iii) $P(H=h \mid D=d) \geq P(D=d \mid H=h) P(H=h)$

According Bayes Rule,

$P(H=h \mid D=d) = \dfrac{P(D=d \mid H=h) P(H=h)}{P(D=d)}$

Since $P(D=d) \leq 1$, thus

$P(H=h \mid D=d) \geq P(D=d \mid H=h) P(H=h)$

3)  (a) Let the jth eigenvector $u_j$, then

(i) $u_j^T A u_j \geq 0$   since $A$ is PSD.

$\Leftrightarrow u_j^T U \Lambda U^T u_j \geq 0$

$\Leftrightarrow u_j^T \left( \sum_{i=1}^{d} \lambda_i u_i u_i^T \right) u_j \geq 0$

$\Leftrightarrow u_j^T \left( \sum_{i=1}^{d} \lambda_i u_i u_i^T u_j \right) \geq 0$

$\Leftrightarrow u_j^T \lambda_j u_j \geq 0$   since $U$ is orthogonal

$\Leftrightarrow \lambda_j \geq 0$   for all $j = 1, 2, \cdots, d$

(ii) if $\lambda_i \geq 0$ for each $i$.

$$A = U \Lambda U^T = \sum_{i=1}^{d} \lambda_i u_i u_i^T$$

let $x \in \mathbb{R}^d$

$$x^T A x = x^T \sum_{i=1}^{d} \lambda_i u_i u_i^T x$$

$$= \sum_{i=1}^{d} \lambda_i \| u_i^T x \|_2^2$$

$$\geq 0 \qquad \text{for all } x$$

Thus $A$ is PSD.

Combine (i) and (ii), we have $A$ is PSD iff $\lambda_i \geq 0$ for each $i$.

(b) (i) If $A$ is PD. By the same arguments before, we have

$u_j^T A u_j > 0$    $u_j$ is the jth eigenvector

$\Leftrightarrow u_j^T \sum_{i=1}^{d} \lambda_i u_i u_i^T u_j > 0$

$\Leftrightarrow u_j^T \lambda_j u_j > 0$

$\Leftrightarrow \lambda_j > 0$   for all $j = 1, 2, \cdots, d$

(ii) If $\lambda_i > 0$ for each $i$.

    let $x \in \mathbb{R}^d$,

$$x^T A x = x^T \sum_{i=1}^{d} \lambda_i u_i u_i^T x$$

$$= \sum_{i=1}^{d} \lambda_i \|u_i^T x\|_2^2$$

        $> 0$          for all $x \neq 0$

Thus $A$ is PD.

Combine (i) and (ii) $A$ is PD iff $\lambda_i > 0$ for each $i$.

4) (a) $f(t\underline{x} + (1-t)\underline{y}) = a^T(t\cdot\underline{x} + (1-t)\cdot\underline{y}) + b$

$\quad\quad\quad = t\cdot a^T\underline{x} + (1-t)\cdot a^T\underline{y} + t\cdot b + (1-t)\cdot b$

$\quad\quad\quad = t(a^T\underline{x} + b) + (1-t)(a^T\underline{y} + b)$

$\quad\quad\quad = t\,f(\underline{x}) + (1-t)\,f(\underline{y})$       ①

Therefore $f(x) = a^T\underline{x} + b$ is convex.

By the same reasoning, we can also have ~~that~~

$\quad -f(t\underline{x} + (1-t)\underline{y}) = -t\,f(\underline{x}) + (1-t)\,f(\underline{y})$   ② .

implies $-f(x)$ is convex, thus $f(x)$ is concave.


From ① ②, the equality holds, so $f(x)$ is not strictly convex.


(b) Suppose there exists more than one global minimizer $x^*$, $y^*$ where $x^* \neq y^*$

we have $\quad f(x^*) = f(y^*)$

Since $f(x)$ is strictly convex on dom $(f)$, we have

$\quad f(t\cdot x^* + (1-t)y^*) < t\,f(x^*) + (1-t)\,f(y^*)$

$\Leftrightarrow\quad f(t(x^* - y^*) + y^*) < f(y^*)$   (or $f(x^*)$ )

$\Rightarrow\quad$ there exists $t(x^* - y^*) + y^* \in$ dom $(f)$ such that

$\quad f(t(x^* - y^*) + y^*) < f(y^*)$

which contradicts that $y^*$ is the global minimizer.

Therefore $f$ has at most one global minimizer.

(c) from (2), we have $f(x) = f(y) + \nabla f^T_{(y)}(x-y) + \frac{1}{2}(x-y)^T \nabla^2 f(y+t(x-y))(x-y)$

let $y=x^*$ and $t=0$, we get

$$f(x) = f(x^*) + \nabla f^T(x^*)(x-x^*) + \frac{1}{2}(x-x^*)^T \nabla^2 f(x^*)(x-x^*)$$

$\Rightarrow f(x)-f(x^*) = \frac{1}{2}(x-x^*)\nabla^2 f(x^*)(x-x^*)$    since $\nabla f(x^*)=0$

$\Rightarrow (x-x^*)\nabla^2 f(x^*)(x-x^*) \geq 0$    since $x^*$ is local minimum.

$\Rightarrow \nabla^2 f(x^*) \succeq 0$

(d) From (2), we get

$$f(x) = f(y) + \nabla f^T(y)(x-y) + \frac{1}{2}(x-y)^T \cdot \nabla^2 f(y+t(x-y))\cdot(x-y) \text{ for all } x,y \in \mathbb{R}^d$$

① Since $f$ is convex, we have

$$f(x) \geq f(y) + \nabla f^T(y)(x-y) \qquad \text{1st order condition.}$$

Therefore, $\frac{1}{2}(x-y)^T \cdot \nabla^2 f(y+t(x-y))(x-y)$

$$= \frac{1}{2}(x-y)^T \cdot \nabla^2 f(tx+(1-t)y)(x-y) \geq 0$$

let $t=1$

$\Rightarrow \nabla^2 f(x) \succeq 0$ for all $x \in \mathbb{R}^d$

② If $\nabla^2 f(x) \succeq 0$ for all $x \in \mathbb{R}^d$, then

$$a^T \cdot \nabla^2 f(x) \cdot a \geq 0$$

$\Rightarrow a^T \cdot \nabla^2 f(y+t(x-y))\cdot a \geq 0$    since $y+t(x-y) \in \mathbb{R}^d$

$\Rightarrow (x-y)^T \nabla^2 f(y+t(x-y))(x-y) \geq 0$

$\Rightarrow f(x) \geq f(y) + \nabla f^T(y)(x-y)$

Therefore $f$ is convex.

Combining ① and ② if $f$ is twice differentiable, then $f$ is convex iff $\nabla^2 f(x)$ is PSD $\forall x \in \mathbb{R}^d$

(e) $\dfrac{\partial f}{\partial x_l \partial x_k}(x) = \dfrac{\partial}{\partial x_l}\left[\dfrac{\partial f(x)}{\partial x_k}\right] = \dfrac{\partial}{\partial x_l}\dfrac{\partial}{\partial x_k}\left[\dfrac{1}{2}x^T A x + b^T x + c\right]$

$\qquad\qquad = \dfrac{\partial}{\partial x_l}\dfrac{\partial}{\partial x_k}\left[\dfrac{1}{2}\sum_{i=1}^{d} x_i \sum_{j=1}^{d} A_{ij} x_j + \sum_{i=1}^{d} b_i x_i + c\right]$

$\qquad\qquad = \dfrac{\partial}{\partial x_l}\left[\sum_{j=1}^{d} A_{kj} x_j + b_k\right]$

$\qquad\qquad = \quad A_{kl} = A_{lk} \qquad$ since $A$ is symmetric.

Therefore $\nabla_x^2 f(x) = A$

According to (d), we have

$\quad f$ is convex if and only if $A$ is positive semi-definite on $\mathbb{R}^d$

Also, by the same argument followed by (d), we conclude

$\quad f$ is strictly convex if and only if $A$ is positive definite on $\mathbb{R}^d$

In [13]:

```python
# Load the Libraries
import math
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

import cv2
```
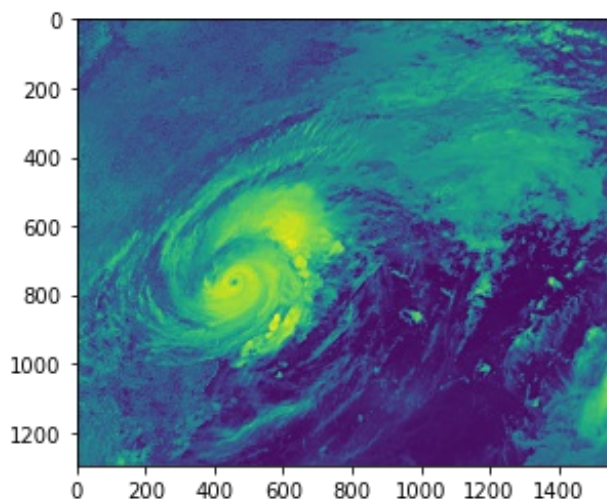
In [6]:

```python
# Load the color image in grayscale
img = cv2.imread('harvey-saturday-goes7am.jpg')
print("The size of the color image is " + str(img.shape))
gray = cv2.cvtColor(img, cv2.COLOR_RGB2GRAY)
print("The size of the grayscale image is " + str(gray.shape))
plt.imshow(gray)
plt.show()
```

```
The size of the color image is (1296, 1548, 3)
The size of the grayscale image is (1296, 1548)
```



In [25]:

```python
# SVD
U, Sigma, V_trans = np.linalg.svd(gray, full_matrices=1)
print(U.shape)
print(Sigma.shape)
print(V_trans.shape)
```

```
(1296, 1296)
(1296,)
(1548, 1548)
```

In [29]:

```python
def F_norm(gray):
    """
    This function calculate the F-norm given a input matrix
    """
    tr_AA_transpose = np.trace(np.matmul(gray,gray.transpose()))
    norm = math.sqrt(tr_AA_transpose)
```

```
    norm = math.sqrt(tr_AA_transpose)
    return norm
print(F_norm(gray))
```
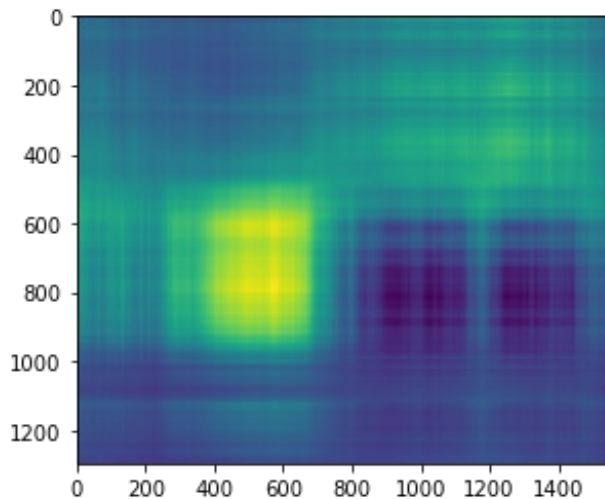
411.85191513455413

In [36]:

```
k = 2
X_bar = 0
for i in range(k):
    X_bar += Sigma[i] * np.outer(U.T[i], V_trans[i])
plt.imshow(X_bar)
plt.show()
```



In [39]:

```
dif = np.subtract(gray,X_bar)
print(F_norm(dif)/F_norm(gray))
```
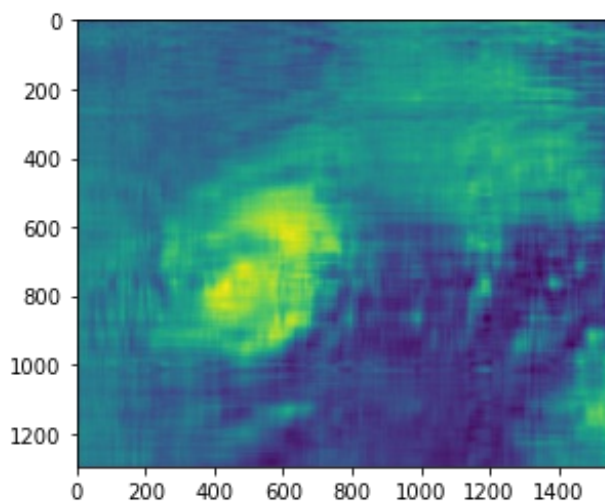
102.4922553316057

In [40]:

```
k = 10
X_bar = 0
for i in range(k):
    X_bar += Sigma[i] * np.outer(U.T[i], V_trans[i])
plt.imshow(X_bar)
plt.show()
```
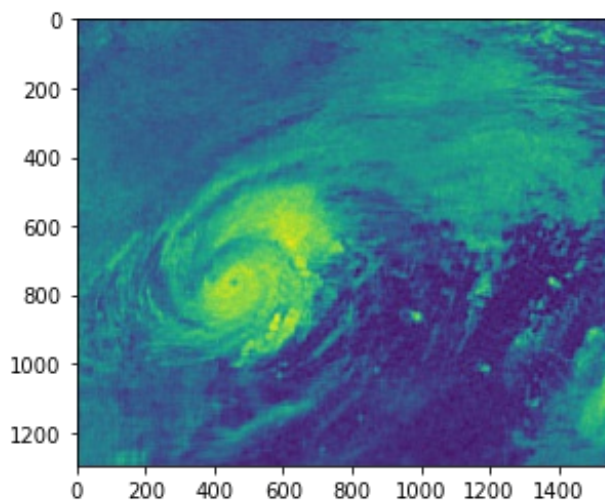
In [41]:

```
dif = np.subtract(gray,X_bar)
print(F_norm(dif)/F_norm(gray))
```

58.93770105042957

In [42]:

```
k = 40
X_bar = 0
for i in range(k):
    X_bar += Sigma[i] * np.outer(U.T[i], V_trans[i])
plt.imshow(X_bar)
plt.show()
```
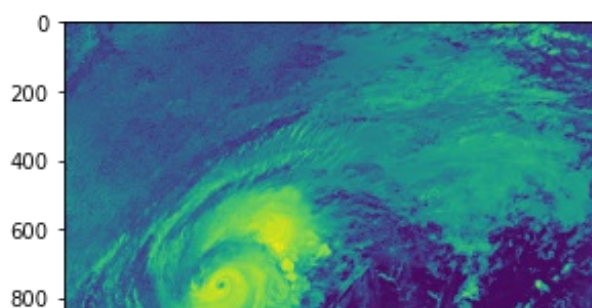


In [43]:

```
dif = np.subtract(gray,X_bar)
print(F_norm(dif)/F_norm(gray))
```
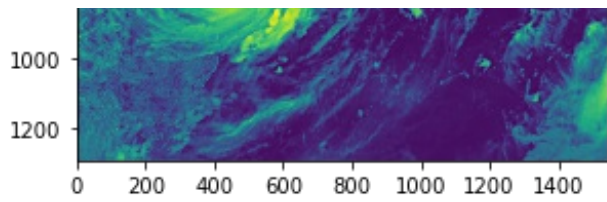
31.300100215056247

In [44]:

```
# this is only for test
k = 1296
X_bar = 0
for i in range(k):
    X_bar += Sigma[i] * np.outer(U.T[i], V_trans[i])
plt.imshow(X_bar)
plt.show()

dif = np.subtract(gray,X_bar)
print(F_norm(dif)/F_norm(gray))
```

(b) According to the fundamental theorem, $(1296 \times k) + (k \times k) + (1548 \times k) = 2844 \times k + k^2$

When k = 2, we need $(1296 \times 2) + (2 \times 2) + (1548 \times 2) = 5692$

When k = 10, we need $(1296 \times 10) + (10 \times 10) + (1548 \times 10) = 28540$

When k = 40, we need $(1296 \times 40) + (40 \times 40) + (1548 \times 40) = 115360$