

# Юстина Иванова

Программист, data scientist

Кейс-стади

Спикер



**Юстина Иванова,**  
Data scientist по  
Компьютерному зрению  
в компании Dataplex  
Выпускница МГТУ им. Баумана  
Msc Artificial Intelligence,  
University of Southampton

# Основные непонятные моменты

1. Матрица ковариаций
2. Свойства матрицы ковариаций
3. Смысл ковариационной матрицы
4. Проецирование данных на вектор
5. Скалярное произведение.
6. Транспонирование матрицы
7. Собственные векторы
8. Собственные значения
9. Спектральное разложение матрицы
10. Разложение Холецкого
11. Декомпозиция матрицы ковариаций
12. Матрица преобразований
13. Теория вероятности
14. Условная вероятность.

## Основные понятия

Дисперсия — среднеквадратичное отклонение от среднего значения (насколько данные разбросаны)

$$\sigma^2(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Ковариация — наличие зависимости между величинами

$$\sigma(x, y) = \frac{1}{n} \sum_{i=1}^n (x - \mu_x)(y - \mu_y)$$

Корреляция — нормированная ковариация, определяет силу зависимости

$$\sigma(x, y) = \frac{Cov(x, y)}{\sqrt{Var(x)} \sqrt{Var(y)}} = \frac{\frac{1}{n} \sum_{i=1}^n (x - \mu_x)(y - \mu_y)}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \mu_y)^2}}$$

## Матрица ковариаций

Матрица корреляций подсчитывается с помощью формул, которые показывают как данные зависят друг от друга в пространстве  $n$  значений (каждый элемент матрицы равен ковариации двух выборок).

$$\Sigma = \begin{bmatrix} \sigma(X_1, X_1) & \sigma(X_1, X_2) & \dots & \sigma(X_1, X_n) \\ \sigma(X_2, X_1) & \sigma(X_2, X_2) & \dots & \sigma(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \sigma(X_n, X_1) & \sigma(X_n, X_2) & \dots & \sigma(X_n, X_n) \end{bmatrix}$$

# Смысл матрицы ковариаций

Большинство решаемых проблем в data science и машинном обучении идет через нахождение матрицы ковариаций:

- Классификационный анализ
  - Регрессионный анализ
- Метод главных компонент
  - Метод Фишера



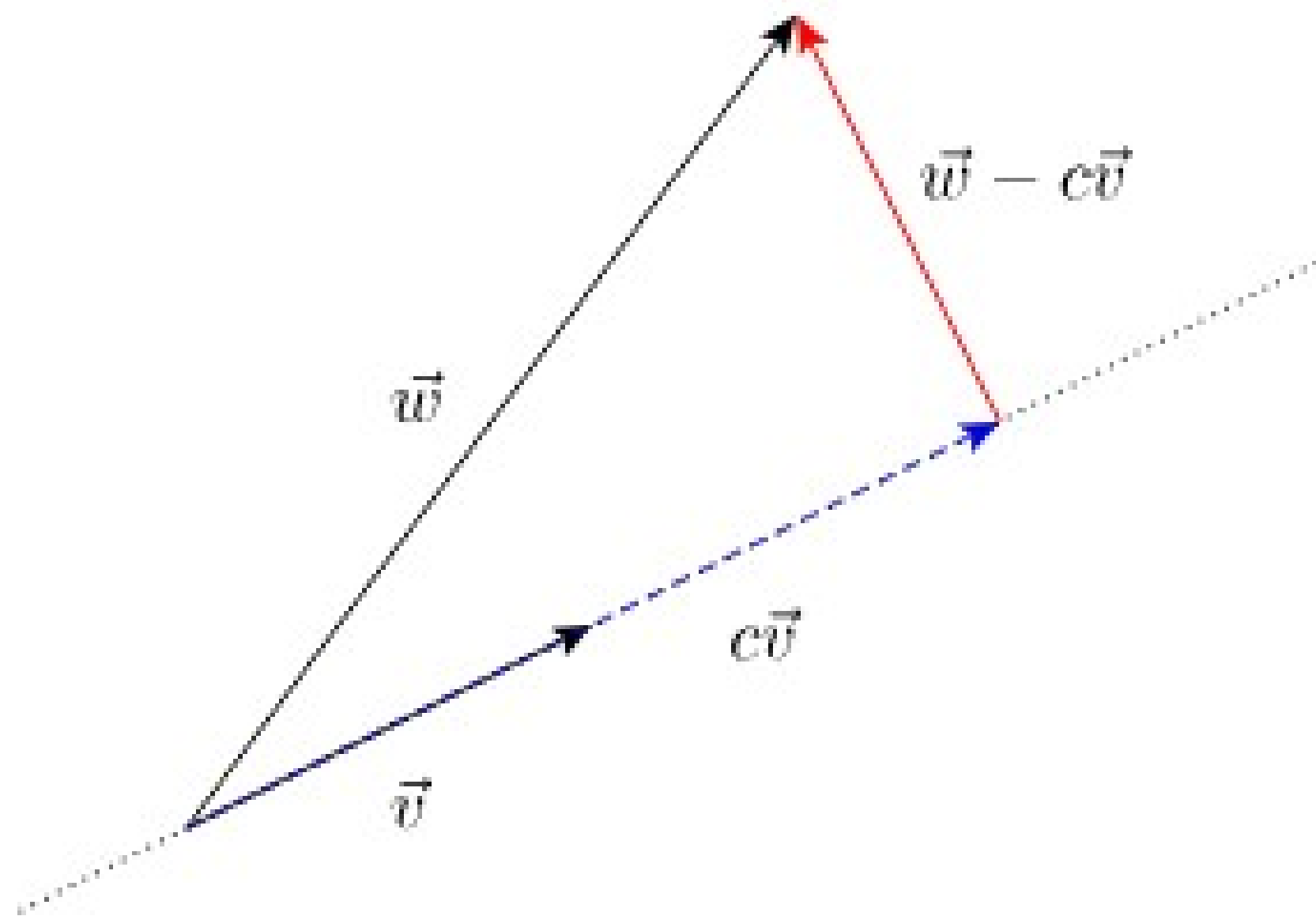
# Классификация

Проблема классификации может решаться такими методами как:

- Дерево решений
- Наивная байесовская классификация
  - Метод наименьших квадратов
  - Метод опорных векторов

# Проецирование данных на вектор

Чтобы посчитать расстояние между точкой и прямой, необходимо знать как проецировать вектор на прямую.

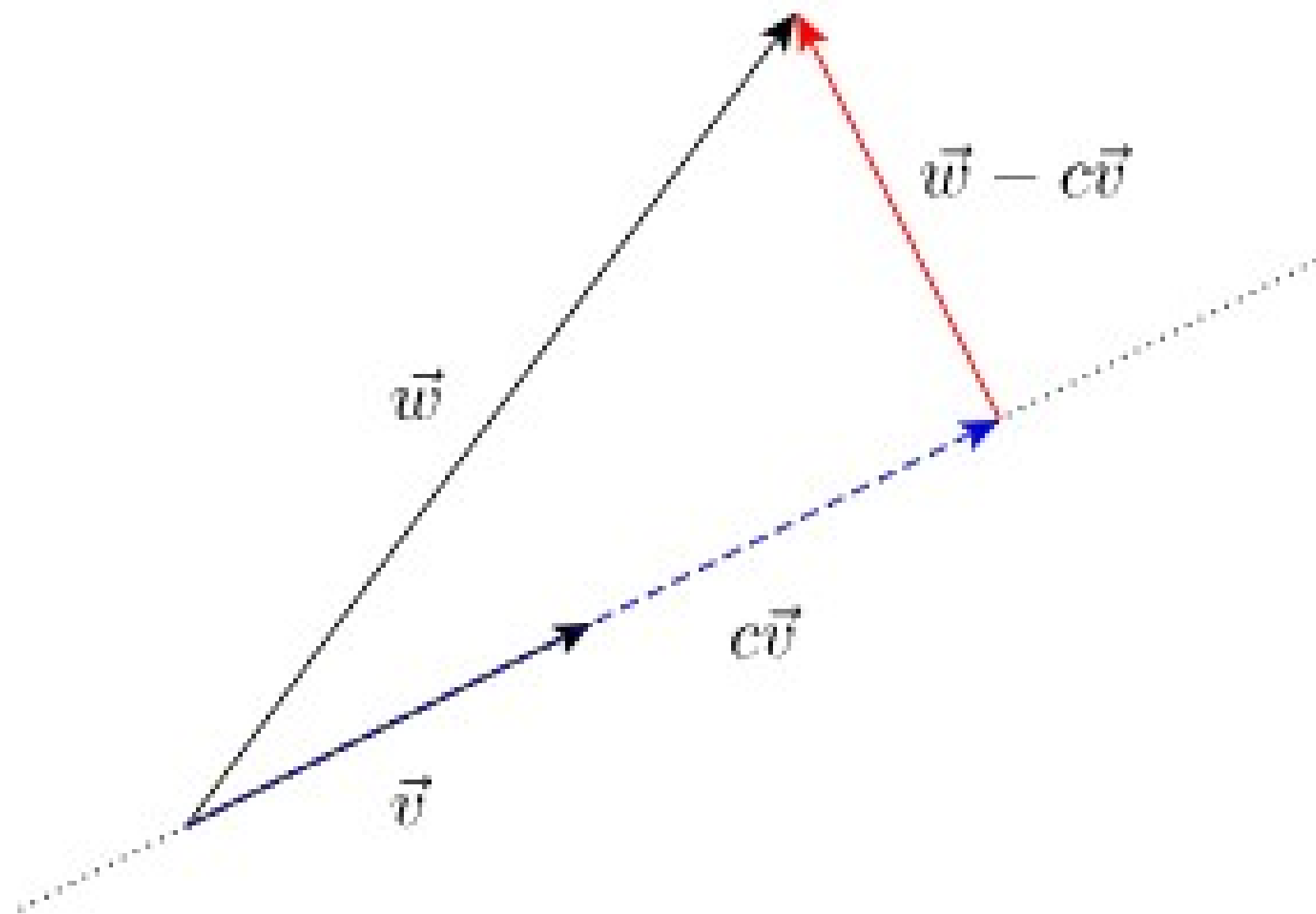


$$\mathbf{cv} = \text{np.dot}(\mathbf{w}, \mathbf{v}) / \text{np.dot}(\mathbf{v}, \mathbf{v}) * \mathbf{v}.$$



# Скалярное произведение

Необходимо для выполнения проецирования данных на вектор.



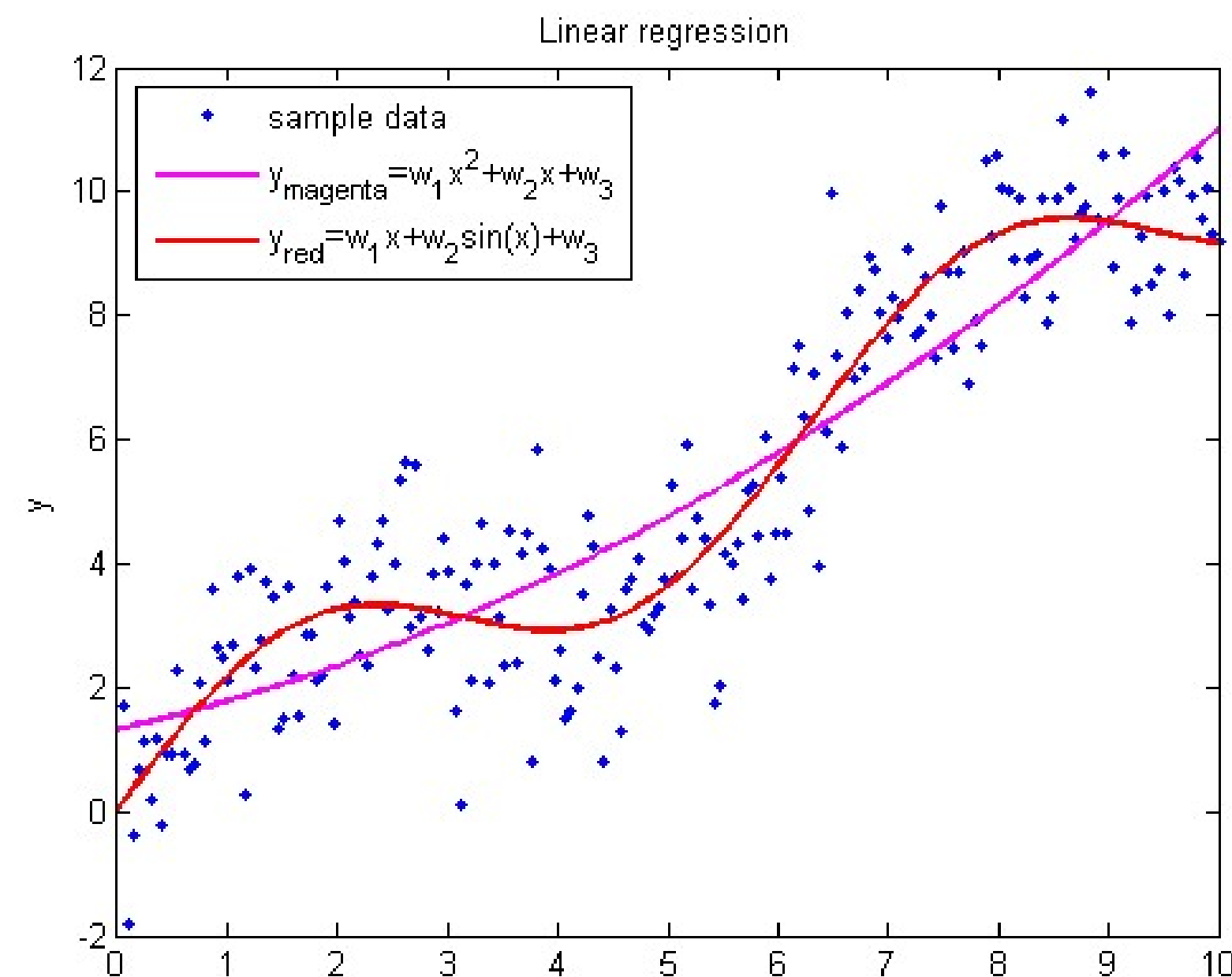
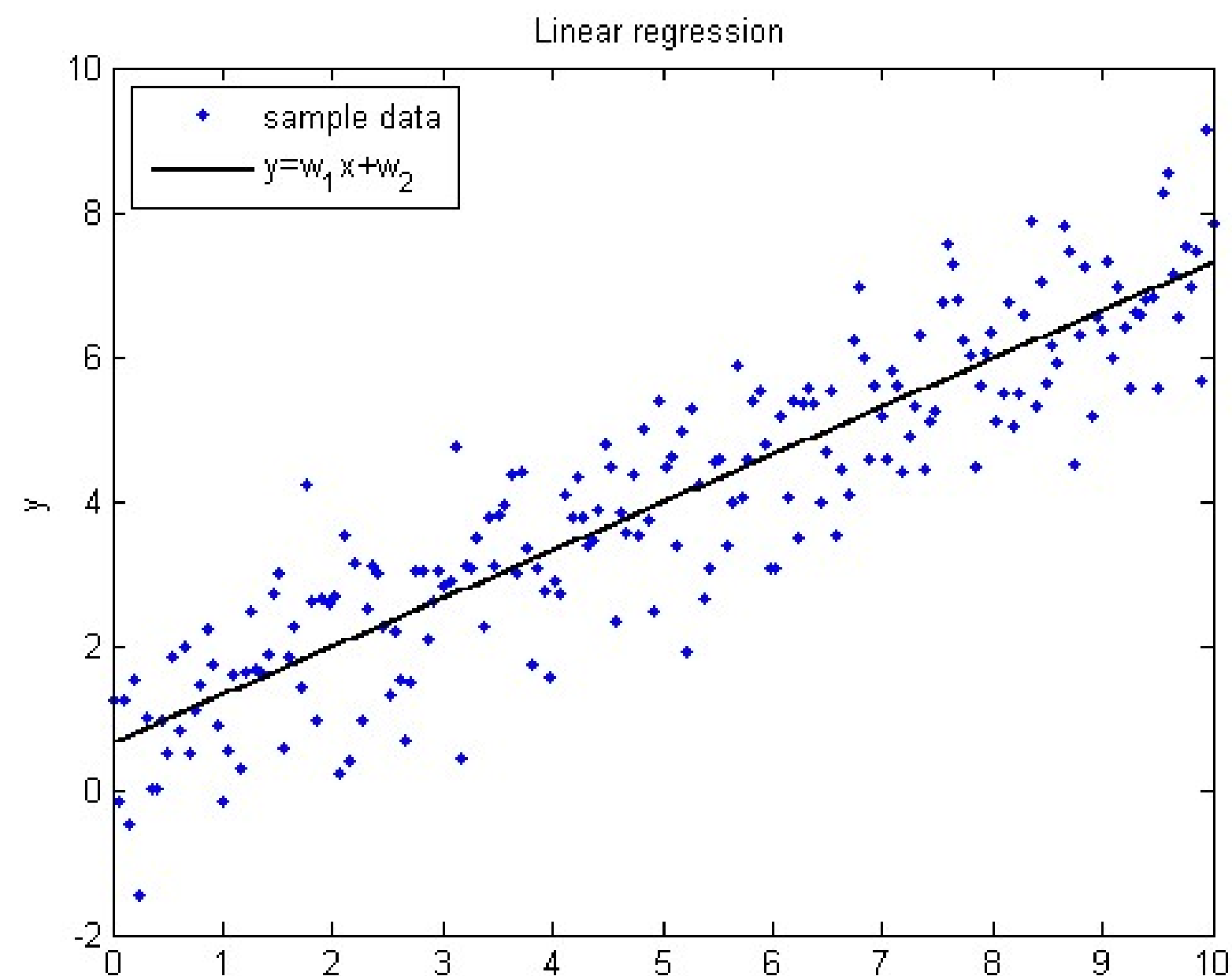
$$\text{np.dot}(w, v) = w_1 * w_2 + v_1 * v_2$$

Где  $w = (w_1, w_2)$

$V = (v_1, v_2)$

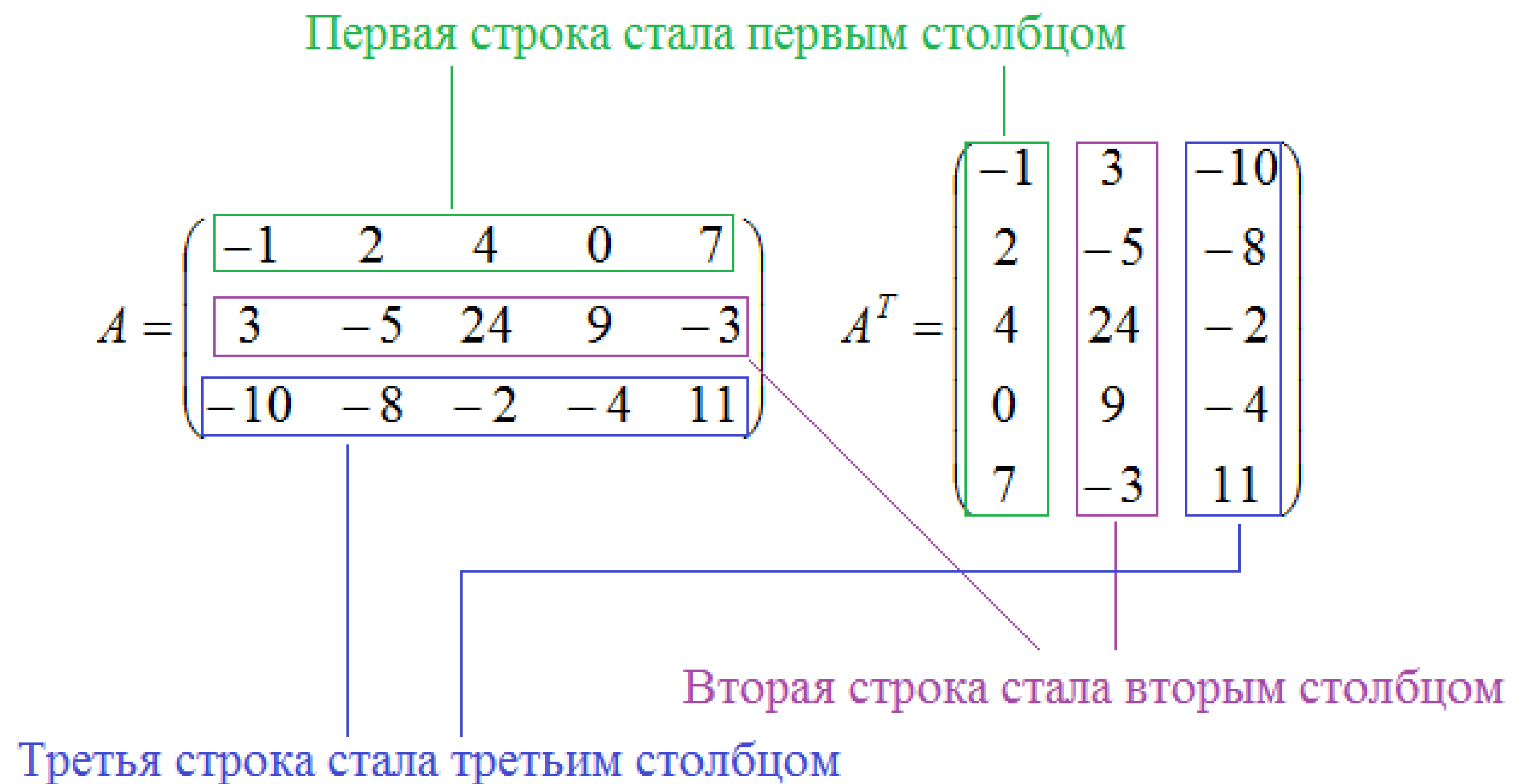
$$\mathbf{cv} = \text{np.dot}(\mathbf{w}, \mathbf{v}) / \text{np.dot}(\mathbf{v}, \mathbf{v}) * \mathbf{v}.$$

# Линейная регрессия



Для заданного пространства данных найти уравнение прямой (или кривой), равномерно разделяющее данные посередине.

# Транспонирование матрицы



Транспонирование матрицы в python:

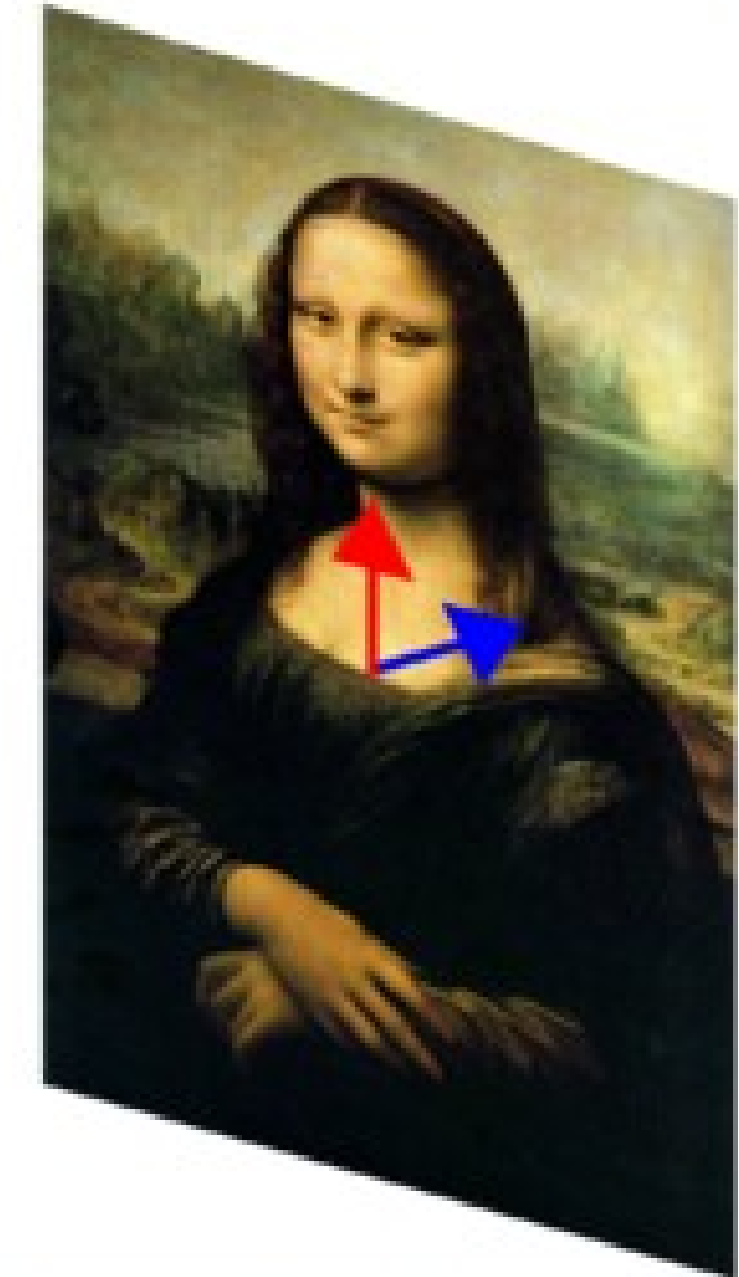
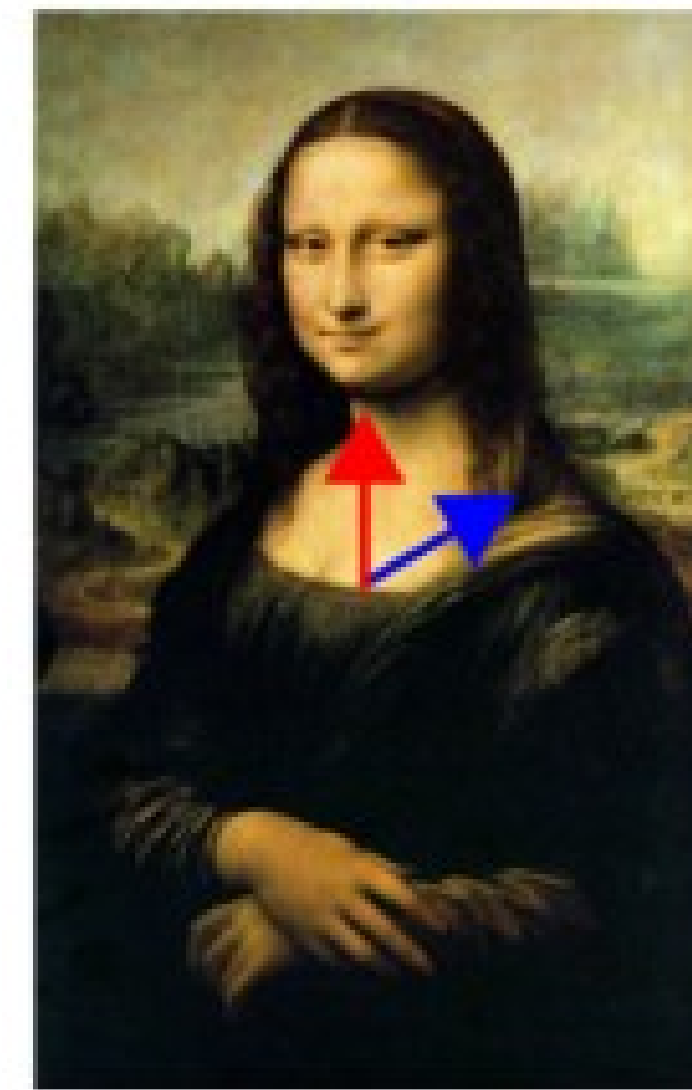
```
import numpy as np  
a = np.array([[0, 1, 2], [4, 5, 6]])  
a = a.transpose()  
или  
a = a.T
```

# Собственные вектора и собственные значения

Считаются через ковариационную матрицу.

Собственным вектором линейного преобразования **A** называется такой ненулевой вектор **x** в **L**, что для некоторого **lambda** в **K**, где **L** — линейное пространство над полем **K**

$$Ax = \lambda x$$



Другая трансформация Джоконды. Синий вектор меняет направление, а красный — нет. Поэтому красный является собственным вектором, а синий — нет. Так как красный вектор ни растянулся, ни сжался, его собственное значение равно, как и на картинке выше, единице. Все векторы, коллинеарные красному, тоже собственные.

# Декомпозиция матрицы ковариаций

Это нахождение собственных векторов и собственных значений.

Scipy и Numpy имеют между собой три различные функции для поиска собственных векторов для заданной квадратной матрицы:

```
numpy.linalg.eig(a)  
scipy.linalg.eig(a) и  
scipy.sparse.linalg.eig(A, k)
```

## Вопросы от участников

Подскажите мне пожалуйста

```
import numpy.linalg as la
```

```
n = 1000
```

`C = [[1,0.98],[0.98,1]]` - Это я так понял матрица ковариации. Почему такие цифры - просто для примера?

```
A = la.cholesky(C)
```

`X = np.random.randn(n,2)` - что это? - генерит 2 случайных числа от одного до 1000? Ваша фраза `X` на самом деле выглядит как `[x, y]` - не понятно

`Y = np.dot(A,X.T)` - Умножаем нашу матрицу на что-то сверху? что такое `X.T`?

`plt.plot(X[:,0], X[:,1], 'r')` - мне непонятно что это такое `X[:,0]`, `X[:,1]` просьба объяснить

```
plt.plot(Y[0,:], Y[1,:], 'b.') (edited)
```



## Ссылки на кейс-стади

[https://github.com/yustiks/statistics\\_in\\_python/blob/master/statistics\\_4.ipynb](https://github.com/yustiks/statistics_in_python/blob/master/statistics_4.ipynb)

## Контакты спикера

yustiks@gmail.com