
Занятие № 11

Feature Selection

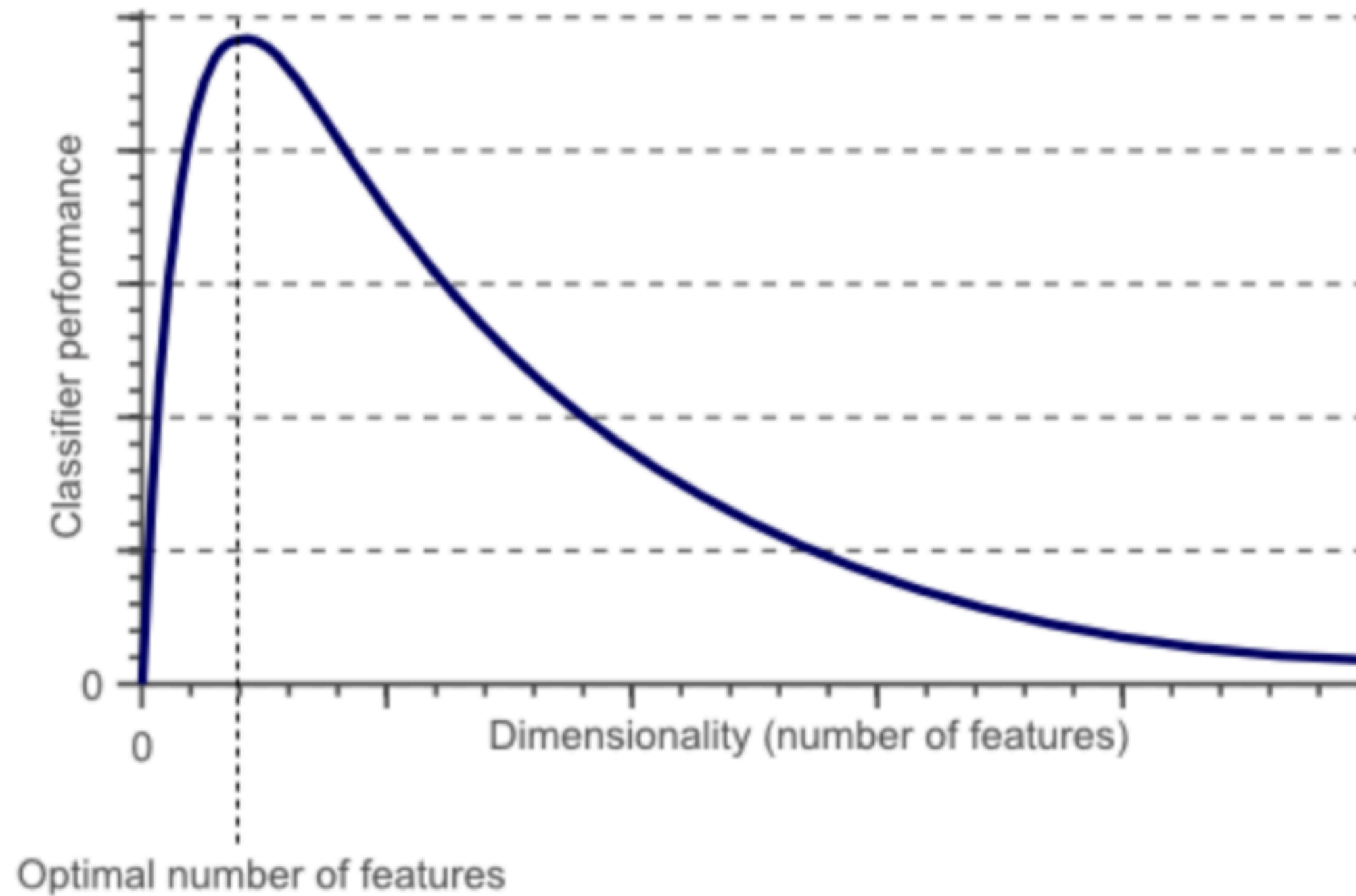


Содержание

- 1 Введение. Зачем всё это?
- 2 Статистика в отборе признаков
- 3 Декомпозиция данных
- 4 Практика.



Введение. Зачем всё это?



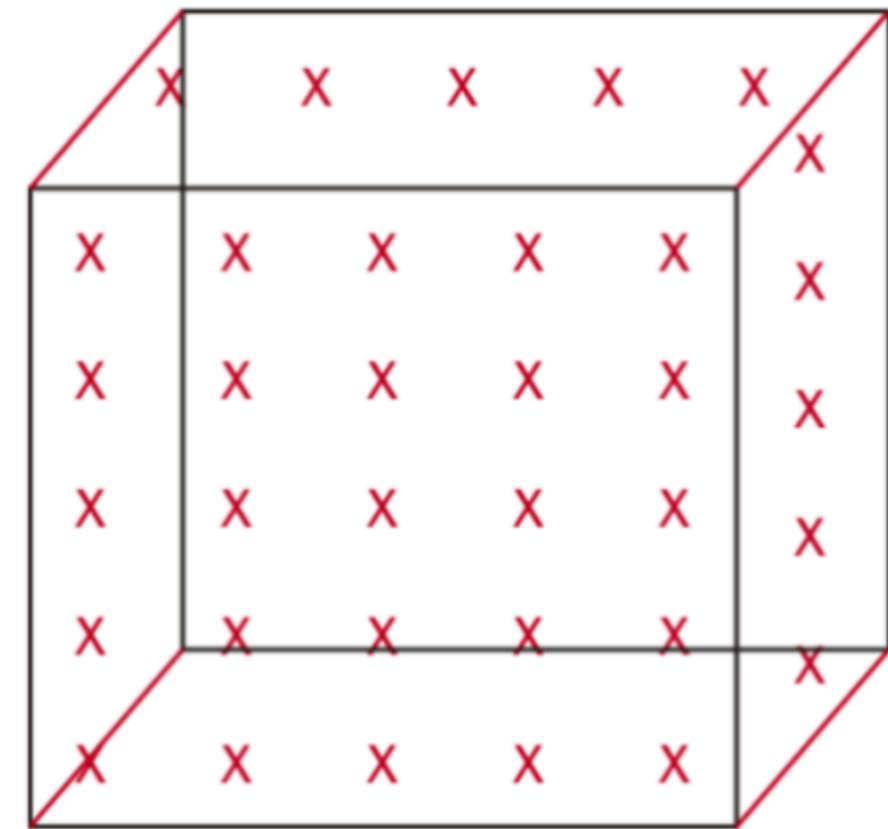
Зачем всё это? Проклятие размерности



Одно измерение - 5 точек



Два измерения - 25 точек



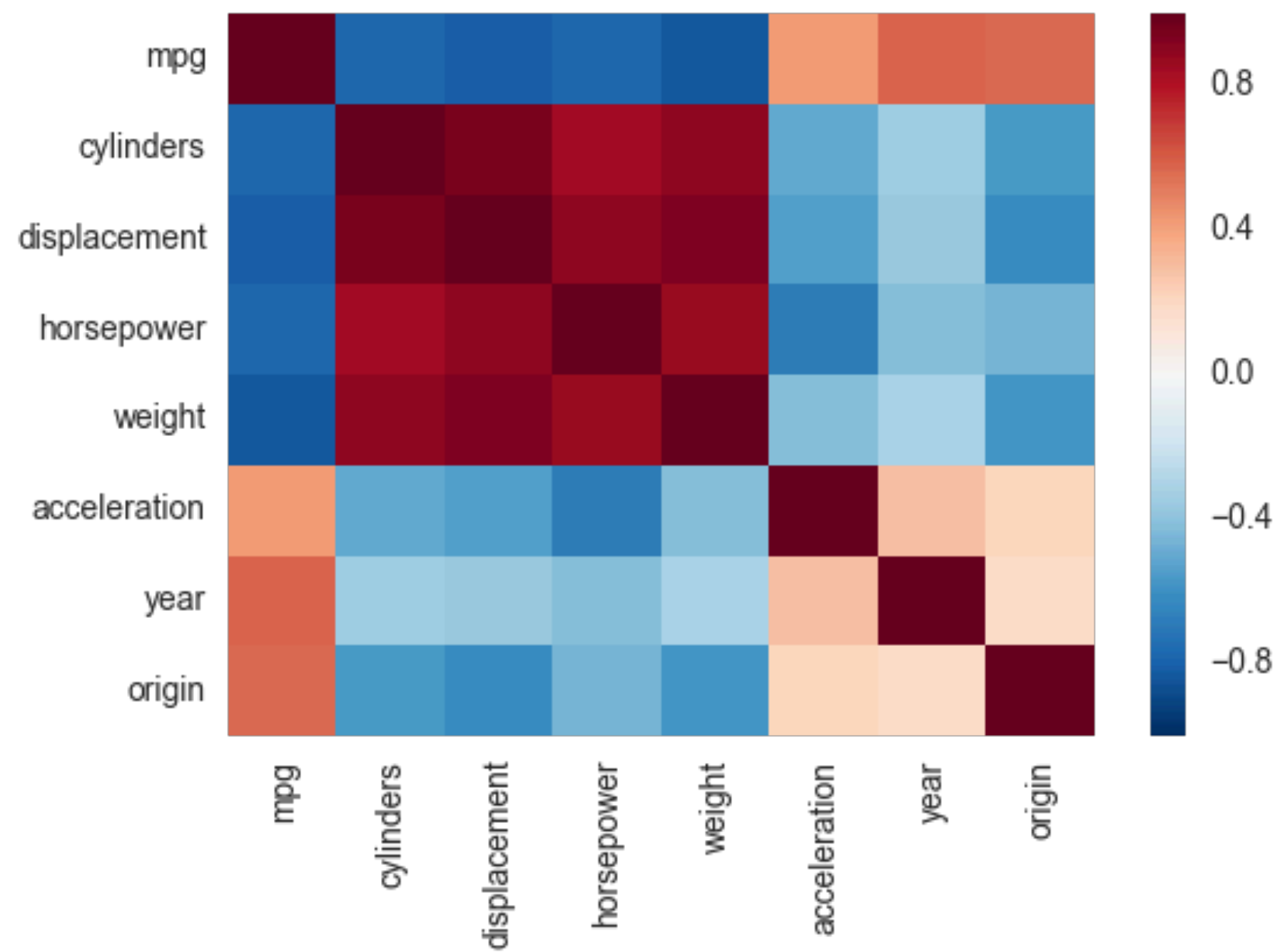
Три измерения - 125 точек



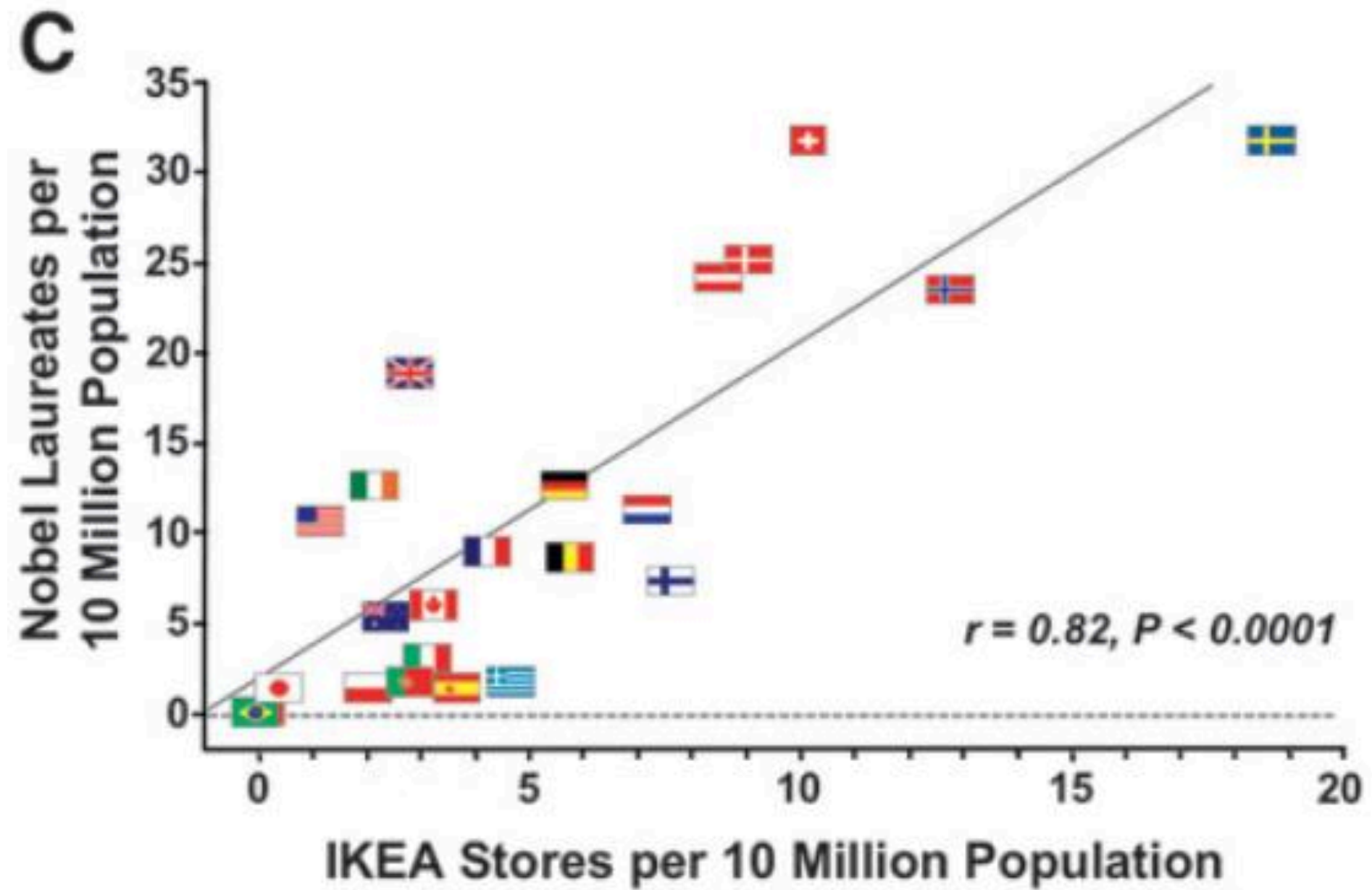
Статистика в отборе признаков



Корреляция



Корреляция



ForexAW.com



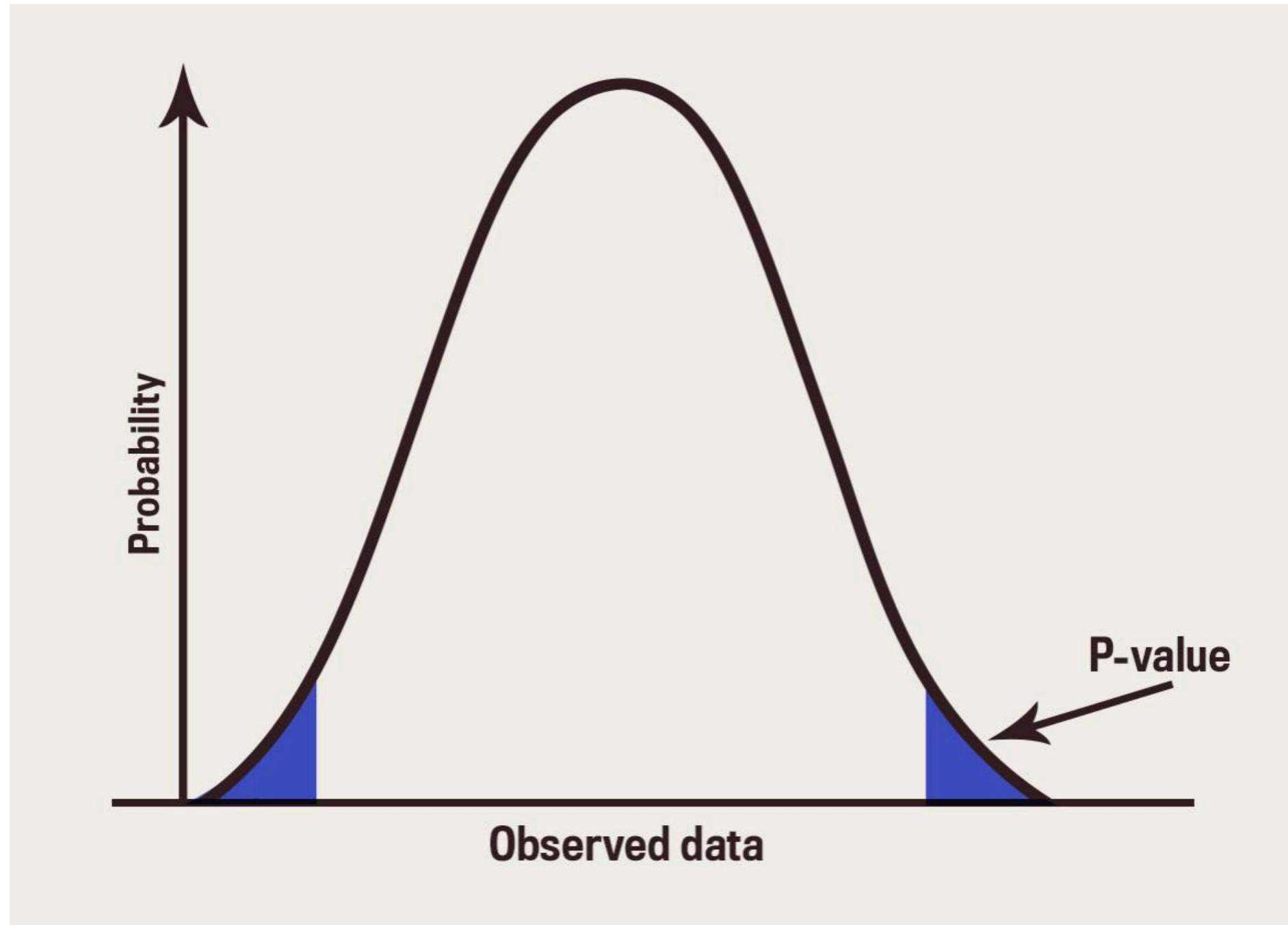
T-статистика

$$t = \frac{\hat{\beta}_i - 0}{SE(\hat{\beta}_i)}$$

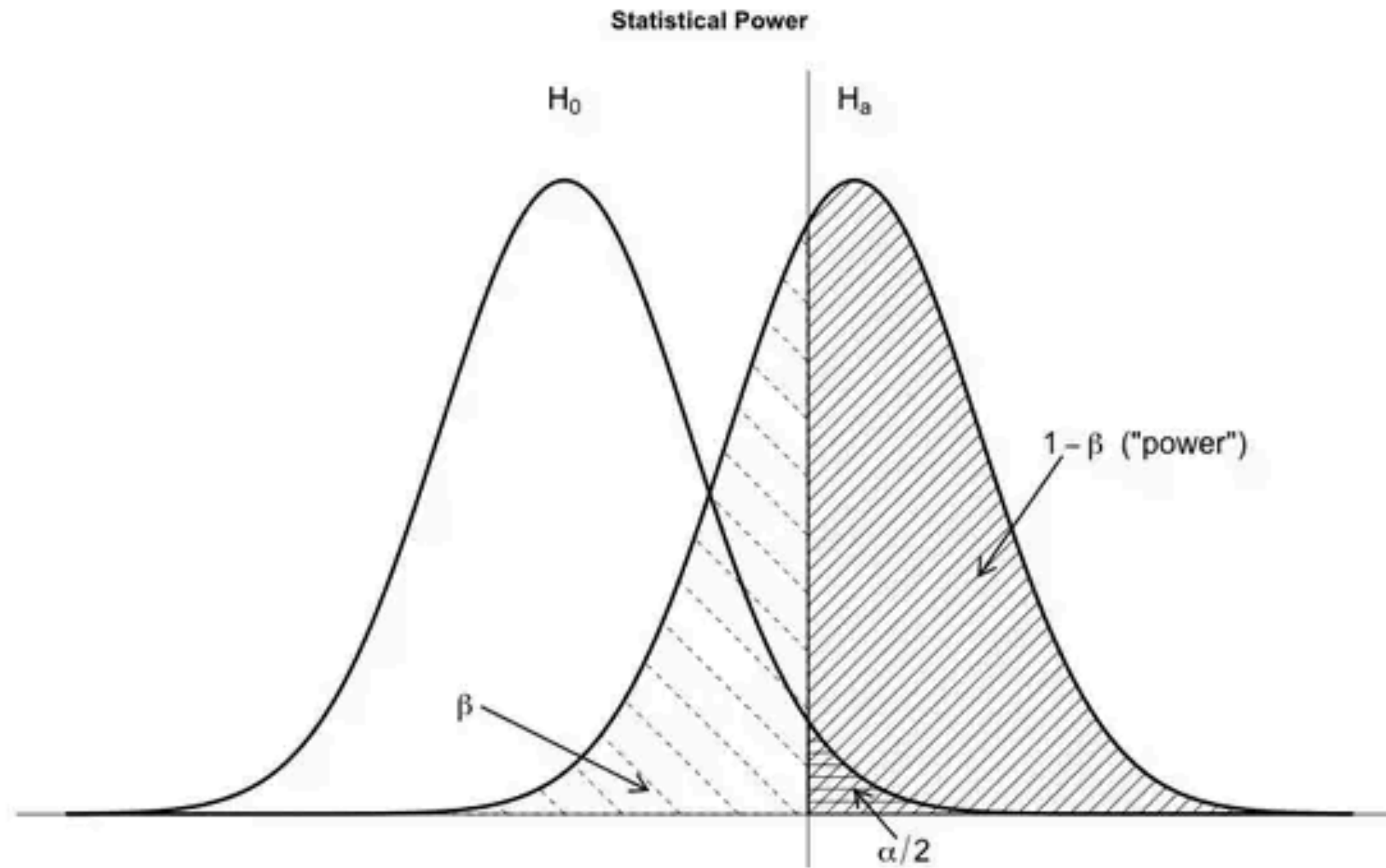
- Если между x_i и y нет зависимости, то t соответствует t -распределению с $n-2$ степенями свободы
- p -value - вероятность того, что при известном распределении наблюдаемое значение $\geq |t|$ (при условии, что $\beta_i = 0$)
- Если p -value достаточно маленький ($< 1\%$), то мы можем отклонить H_0



P-value



P-value

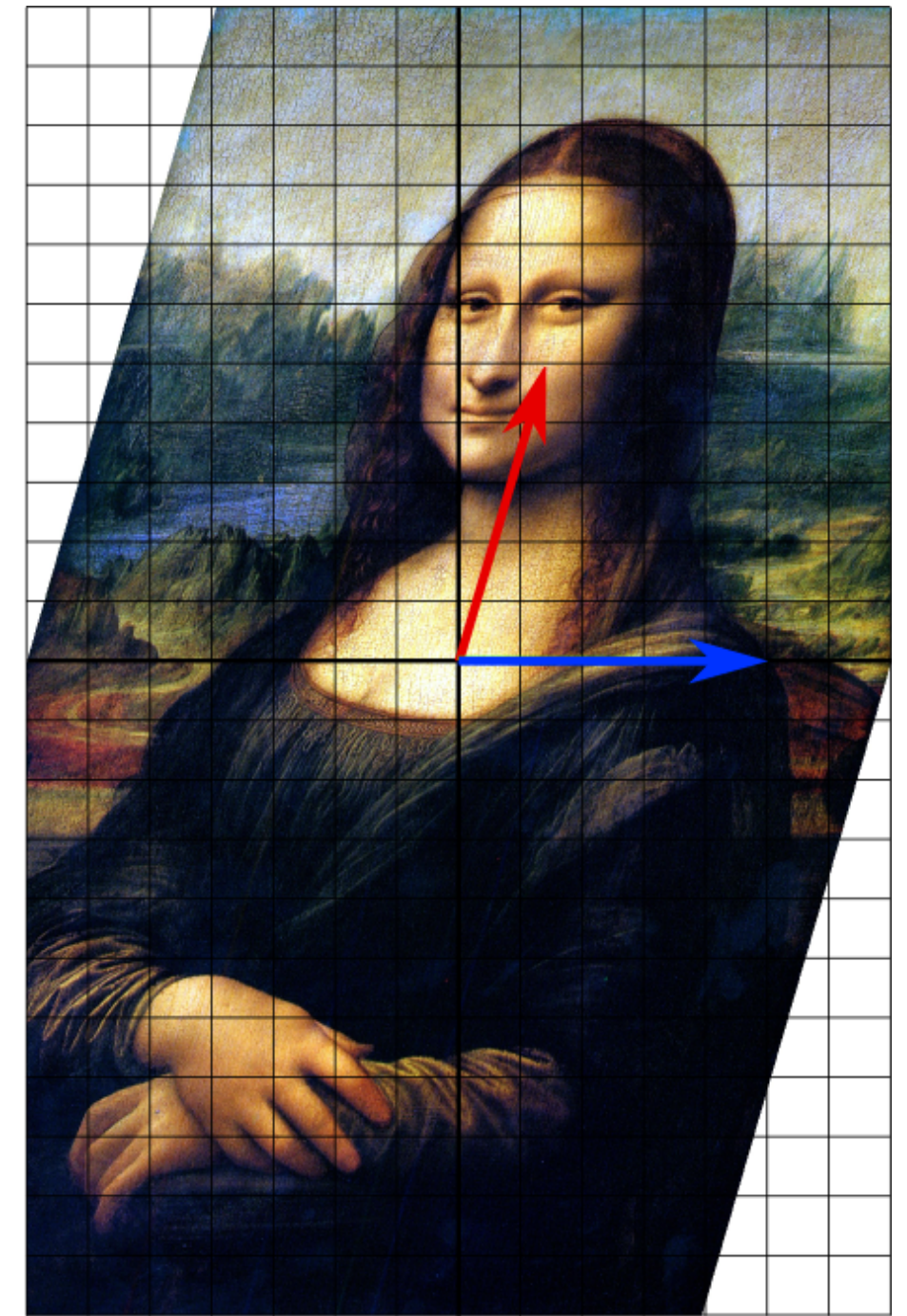
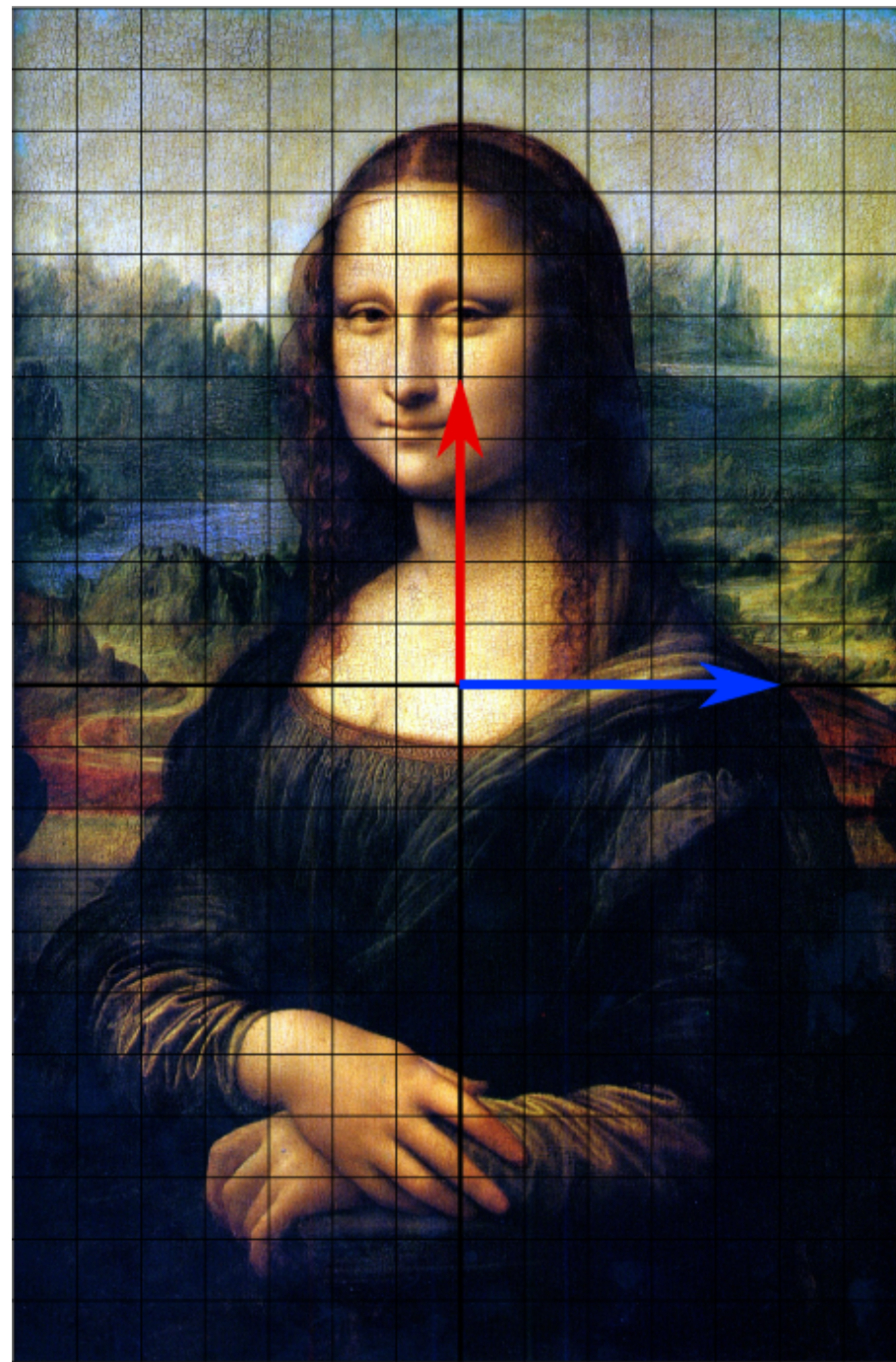


Декомпозиция данных



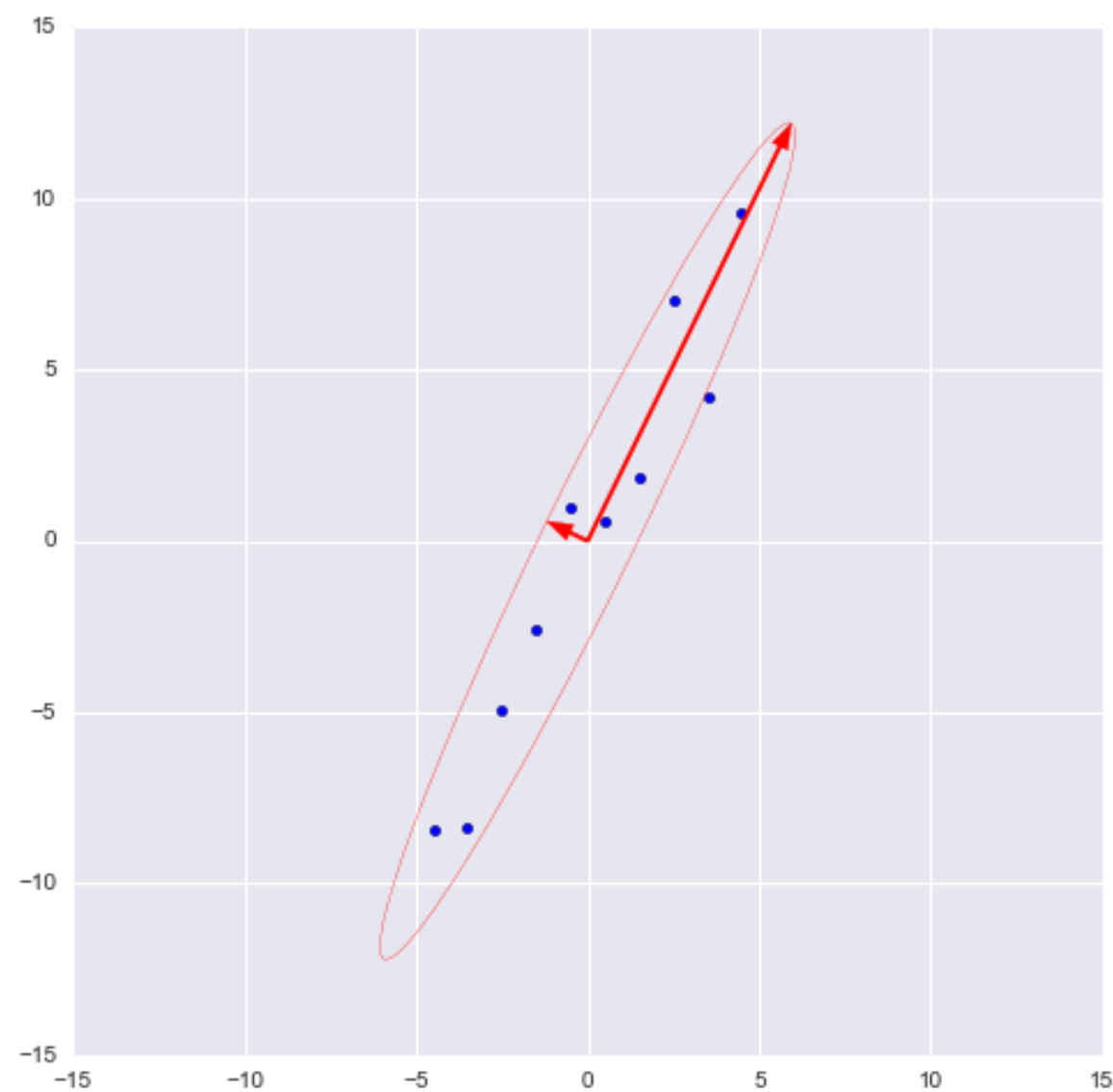
Собственный вектор

$$M\vec{x} = \lambda\vec{x}$$



РСА

Зачем он нужен? Он уменьшает размерность! 😊



PCA

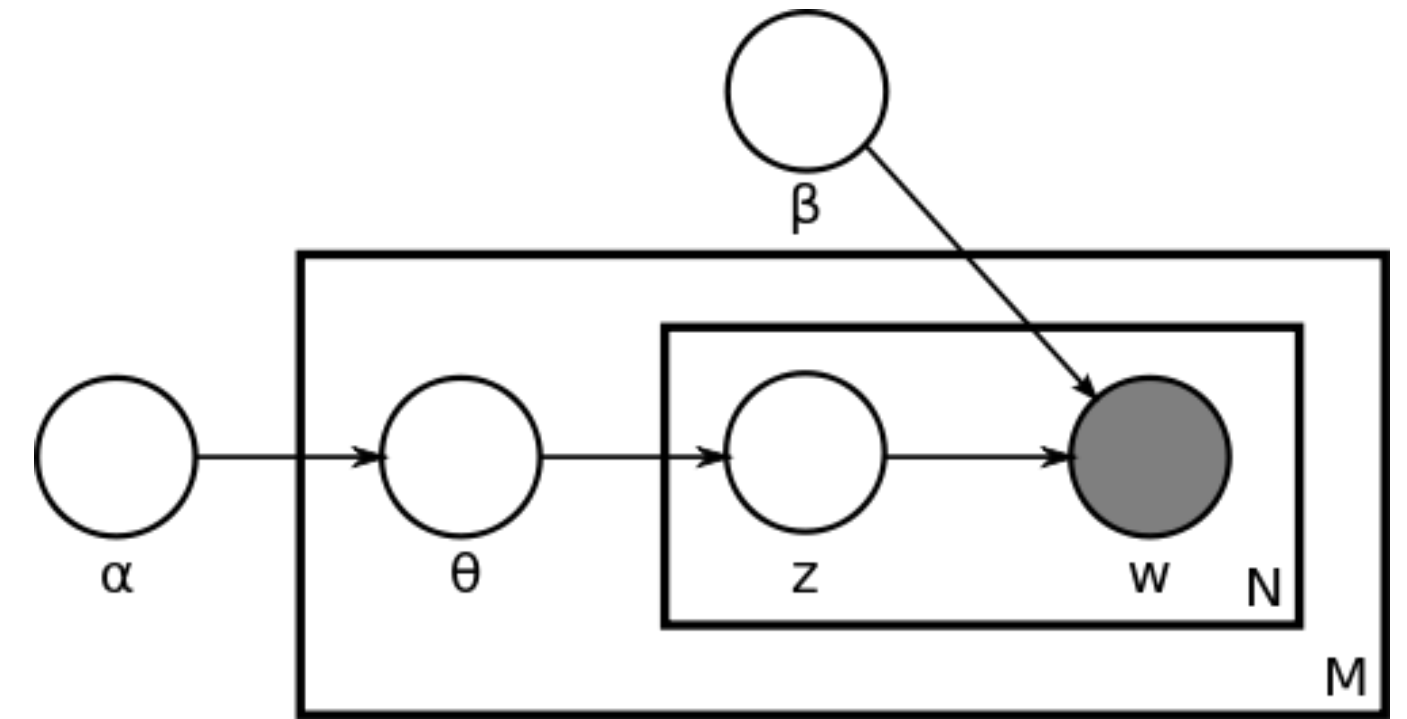
$$\text{Cov}(X_i, X_j) = E\left[(X_i - E(X_i)) \cdot (X_j - E(X_j))\right] = E(X_i X_j) - E(X_i) \cdot E(X_j)$$

$$\begin{aligned}\text{Var}(X^*) &= \Sigma^* = E(X^* \cdot X^{*T}) = E\left((\vec{v}^T X) \cdot (\vec{v}^T X)^T\right) = \\ &= E(\vec{v}^T X \cdot X^T \vec{v}) = \vec{v}^T E(X \cdot X^T) \vec{v} = \vec{v}^T \Sigma \vec{v}\end{aligned}$$



LDA

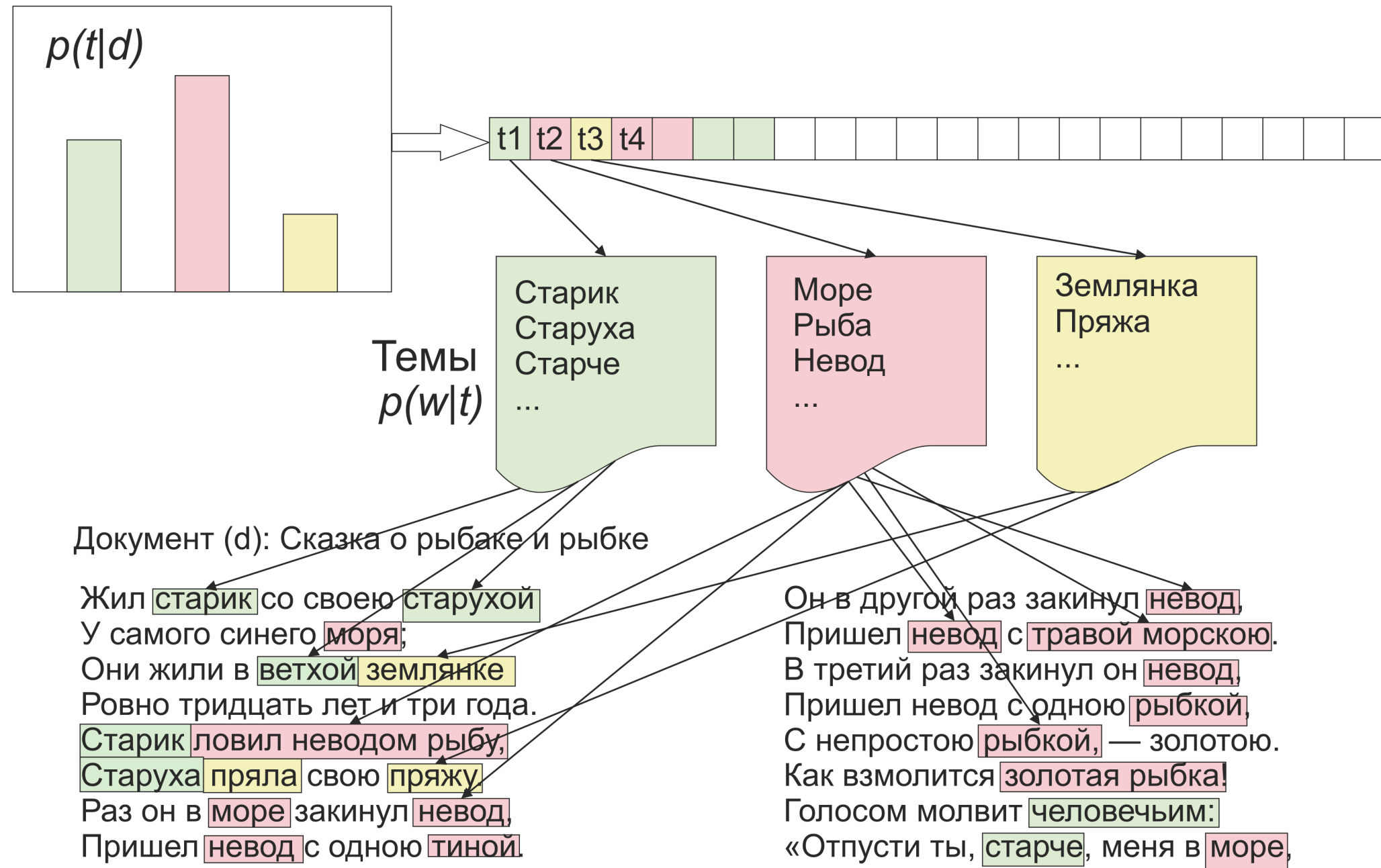
это иерархическая байесовская модель, состоящая из двух уровней:
на первом уровне – смесь, компоненты которой соответствуют «темам»;
на втором уровне – мультиномиальная переменная с априорным распределением Дирихле, которое задаёт «распределение тем» в документе.



$$p(\theta, z, w, N \mid \alpha, \beta) = p(N \mid \xi) p(\theta \mid \alpha) \prod_{n=1}^N p(z_n \mid \theta) p(w_n \mid z_n, \beta).$$



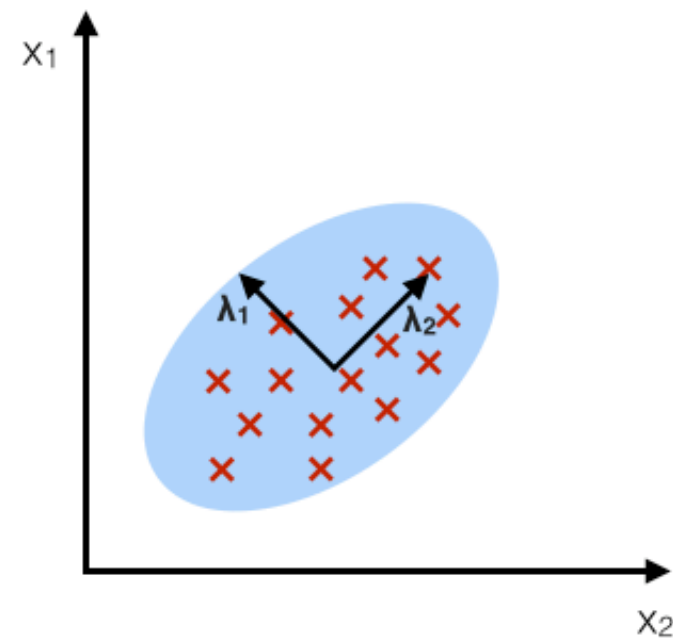
LDA



Сравнение LDA и PCA

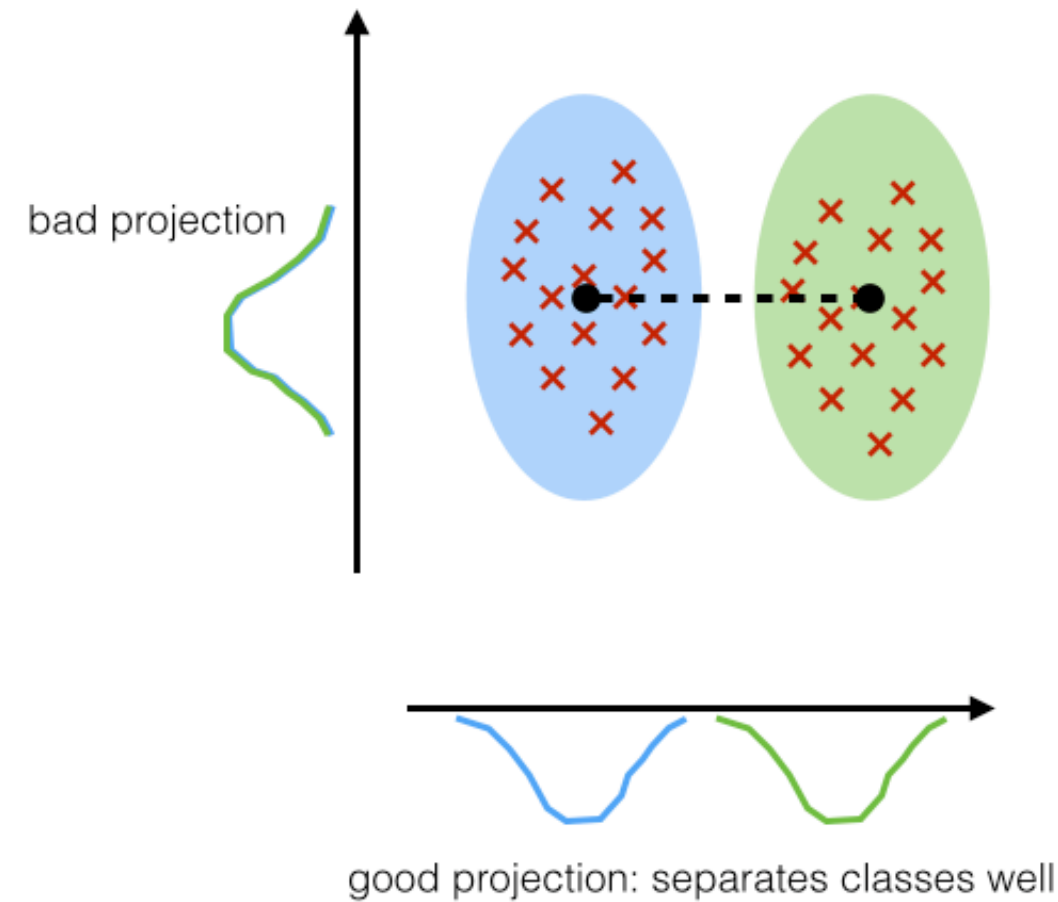
PCA:

component axes that maximize the variance

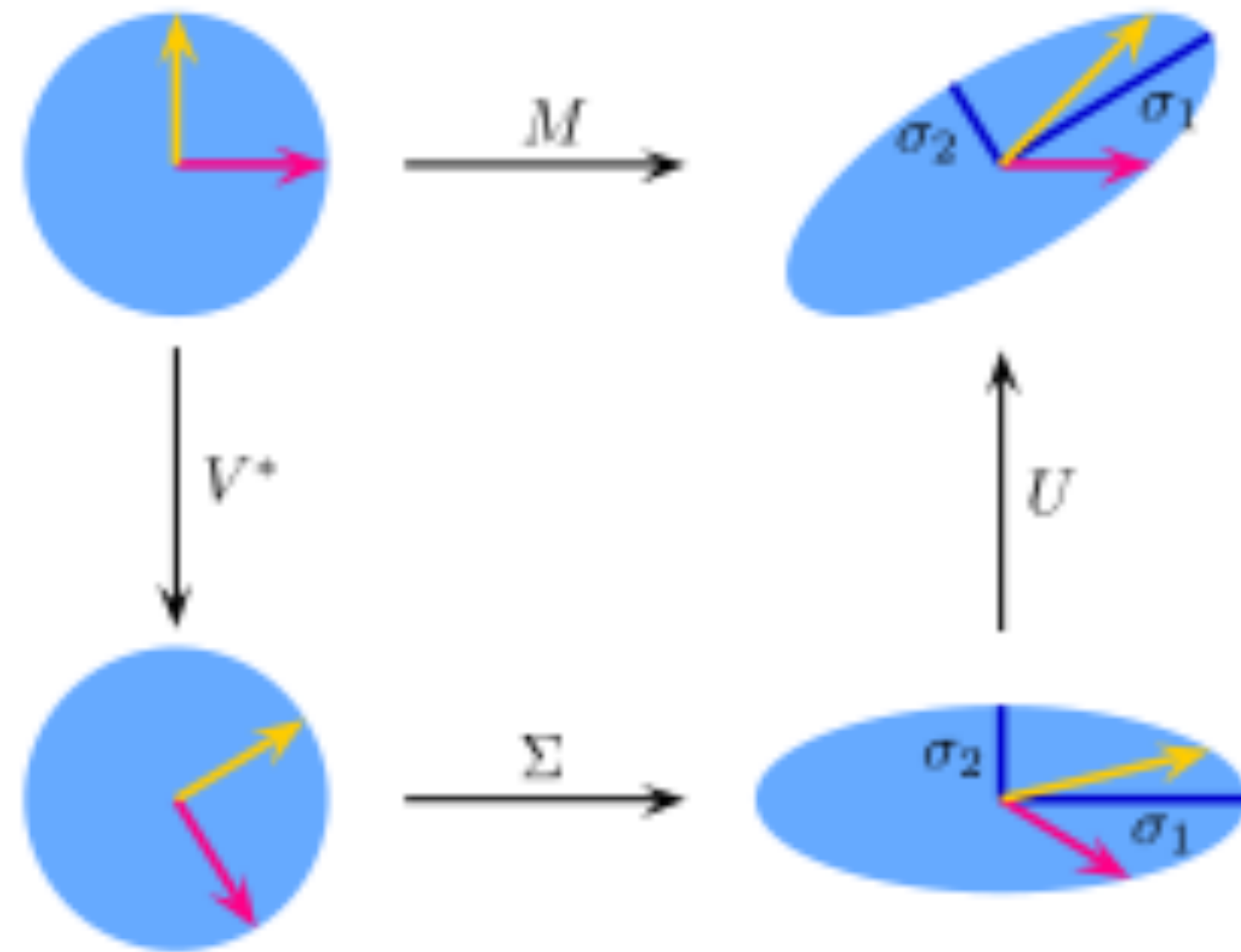


LDA:

maximizing the component axes for class-separation



SVD



$$M = U \cdot \Sigma \cdot V^*$$



SVD

The diagram illustrates the Singular Value Decomposition (SVD) of a matrix A . Matrix A is shown as a single pink rectangle with dimensions $n \times d$. It is equal to the product of three matrices: U , Σ , and V^T . Matrix U is represented by a pink rectangle \hat{U} of size $n \times r$ and a light blue rectangle of size $n \times n$. Matrix Σ is represented by a pink rectangle $\hat{\Sigma}$ of size $r \times r$ and a light blue rectangle of size $n \times d$. Matrix V^T is represented by a pink rectangle \hat{V}^T of size $r \times d$ and a light blue rectangle of size $d \times d$.

$$\begin{matrix} \boxed{\begin{matrix} A \\ n \times d \end{matrix}} = \boxed{\begin{matrix} \hat{U} \\ n \times r \end{matrix}} \boxed{\begin{matrix} \hat{\Sigma} \\ r \times r \end{matrix}} \boxed{\begin{matrix} \hat{V}^T \\ r \times d \end{matrix}} \\ \begin{matrix} U \\ n \times n \end{matrix} \quad \begin{matrix} \Sigma \\ n \times d \end{matrix} \quad \begin{matrix} V^T \\ d \times d \end{matrix} \end{matrix}$$



ПРАКТИКА



Спасибо за внимание!

Сапрыкин Артур
Data Scientist



fb.com/asaprykin92



asaprykin92@gmail.com

