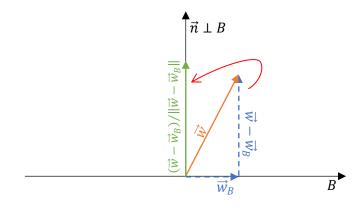
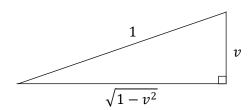
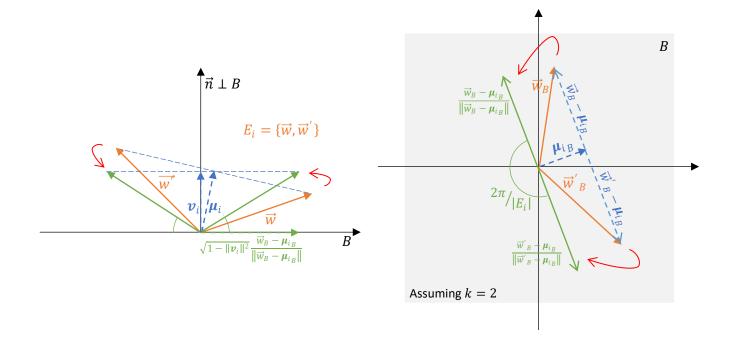
Hard Debiasing

$$\overrightarrow{w} \coloneqq (\overrightarrow{w} - \overrightarrow{w}_B) / \|\overrightarrow{w} - \overrightarrow{w}_B\| \quad , \; \forall \; \; w \in N$$



$$\overrightarrow{w} \coloneqq v_i + \sqrt{1 - \|v_i\|^2} \frac{\overrightarrow{w}_B - \mu_{i_B}}{\|\overrightarrow{w}_B - \mu_{i_B}\|}$$





Identify Bias Subspace

a) $B := \text{first } k \text{ rows of } \mathbb{SVD}(C)$

Let the data matrix \mathbf{X} be of $n \times p$ size, where n is the number of samples and p is the number of variables. Let us assume that it is *centered*, i.e. column means have been subtracted and are now equal to zero.

Then the $p \times p$ covariance matrix \mathbf{C} is given by $\mathbf{C} = \mathbf{X}^{\top} \mathbf{X}/(n-1)$. It is a symmetric matrix and so it can be diagonalized:

$$C = VLV^{T}$$
,

where ${\bf V}$ is a matrix of eigenvectors (each column is an eigenvector) and ${\bf L}$ is a diagonal matrix with eigenvalues λ_i in the decreasing order on the diagonal. The eigenvectors are called *principal axes* or *principal directions* of the data. Projections of the data on the principal axes are called *principal components*, also known as *PC scores*; these can be seen as new, transformed, variables. The j-th principal component is given by j-th column of ${\bf XV}$. The coordinates of the i-th data point in the new PC space are given by the i-th row of ${\bf XV}$.

If we now perform singular value decomposition of \mathbf{X} , we obtain a decomposition

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathsf{T}},$$

where ${f U}$ is a unitary matrix and ${f S}$ is the diagonal matrix of singular values s_i . From here one can easily see that

$$\mathbf{C} = \mathbf{V}\mathbf{S}\mathbf{U}^{\top}\mathbf{U}\mathbf{S}\mathbf{V}^{\top}/(n-1) = \mathbf{V}\frac{\mathbf{S}^2}{n-1}\mathbf{V}^{\top},$$

meaning that right singular vectors ${\bf V}$ are principal directions and that singular values are related to the eigenvalues of covariance matrix via $\lambda_i=s_i^2/(n-1)$. Principal components are given by ${\bf X}{\bf V}={\bf U}{\bf S}{\bf V}^{\top}{\bf V}={\bf U}{\bf S}$.

Reference: https://bit.ly/2IJeAap

b)
$$C \coloneqq \sum_{i=1}^n \sum_{w \in D_i} (\overrightarrow{w} - \mu_i)^T (\overrightarrow{w} - \mu_i) / |D_i|$$