

### Problem Set 3

For the Exercise Sessions on Oct 10 and Oct 17

Last name	First name	SCIPER Nr	Points

#### Problem 1: Some review problems on linear algebra

(a) (*Frobenius norm*) Prove that  $\|A\|_F^2 = \text{trace}(A^H A)$ .

(b) (*Singular Value Decomposition*) Let  $\sigma_i(A)$  denote the  $i^{\text{th}}$  singular value of an  $m \times n$  matrix  $A$ . Prove that  $\|A\|_F^2 = \sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A)$

(c) (*Projection Matrices*) Consider a set of  $k$  orthonormal vectors in  $\mathbb{C}^n$ , denoted by  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ . The projection matrix (that projects an arbitrary vector into the subspace spanned by these orthonormal vectors) is given by

$$P = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^H. \quad (1)$$

- Prove that this matrix is *Hermitian*, i.e.,  $P^H = P$ .
- Prove that this matrix is *idempotent*, i.e.,  $P^2 = P$ . (In words, projecting twice into the same subspace is the same as projecting only once.)
- Prove that  $\text{trace}(P) = k$ , i.e., equal to the dimension of the subspace.
- Prove that the diagonal entries of  $P$  must be real-valued and non-negative. Then, prove that the diagonal entries of  $P$  cannot be larger than 1 (this is a little more tricky).

**Solution 1.** (a) Let  $A$  be an  $m \times n$  matrix and denote by  $a_{ij}$  the entry of  $A$  at row  $i$  and column  $j$ . Hence, we have

$$A^H A = \begin{bmatrix} \sum_{i=1}^m |a_{i1}|^2 & \sum_{i=1}^m a_{i1}^* a_{i2} & \dots & \sum_{i=1}^m a_{i1}^* a_{in} \\ \sum_{i=1}^m a_{i2}^* a_{i1} & \sum_{i=1}^m |a_{i2}|^2 & \dots & \sum_{i=1}^m a_{i2}^* a_{in} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^m a_{in}^* a_{i1} & \sum_{i=1}^m a_{in}^* a_{i2} & \dots & \sum_{i=1}^m |a_{in}|^2 \end{bmatrix} \quad (2)$$

$$\text{trace}(A^H A) = \sum_{j=1}^n \sum_{i=1}^m |a_{ij}|^2 = \left( \sqrt{\sum_{j=1}^n \sum_{i=1}^m |a_{ij}|^2} \right)^2 = \|A\|_F^2 \quad (3)$$

(b) Let  $U \Sigma V^H$  be the singular value decomposition of  $A$ . Then  $U$  is an  $m \times m$  unitary matrix (i.e.  $U^H U = I$ ),  $\Sigma$  is a diagonal  $m \times n$  matrix ( $\Sigma_{ii} = \sigma_i(A)$ ), and  $V$  is an  $n \times n$  unitary matrix (i.e.

$V^H V = I$ ). using part (a) we have

$$\|A\|_F^2 = \text{trace}(A^H A) = \text{trace}((U\Sigma V^H)^H U\Sigma V^H) \quad (4)$$

$$= \text{trace}(V\Sigma^H U^H U\Sigma V^H) \quad (5)$$

$$= \text{trace}(\Sigma^H \Sigma) \quad (6)$$

$$= \sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A) \quad (7)$$

$$(8)$$

where we use the property of unitary matrix and the cyclic property of trace.

(c) For the first bullet item,  $P^H = \left(\sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^H\right)^H = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^H = P$ .

For the second bullet item,  $P^2 = \sum_{i=1}^k \sum_{j=1}^k \mathbf{u}_i \mathbf{u}_i^H \mathbf{u}_j \mathbf{u}_j^H$ . Since  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  are orthonormal vectors, we have

$$\mathbf{u}_i^H \mathbf{u}_j = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases} \quad (9)$$

Hence, all the terms with  $i = j$  survive.  $P^2 = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^H = P$

For the third bullet item,  $\text{trace}(P) = \text{trace}(\sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^H) = \sum_{i=1}^k \text{trace}(\mathbf{u}_i \mathbf{u}_i^H) = \sum_{i=1}^k \text{trace}(\mathbf{u}_i^H \mathbf{u}_i) = k$ .

For the last bullet item, note that we can express the diagonal elements in the following form:

$$P_{11} = \mathbf{e}_1^H P \mathbf{e}_1, \quad (10)$$

where  $\mathbf{e}_1$  is the vector  $(1, 0, 0, \dots, 0)^T$ . Moreover, we know that  $P = P^H = P^2 = P^H P$ . Therefore,

$$P_{11} = \mathbf{e}_1^H P \mathbf{e}_1 = \mathbf{e}_1^H P^H P \mathbf{e}_1 = \|P \mathbf{e}_1\|^2, \quad (11)$$

which establishes that  $P_{11}$  is real-valued and non-negative. The tricky part is now to show that it is also upper bounded by 1.

An elegant proof of this fact follows by considering the matrix  $Q = I - P$ . Clearly, if we can show that the diagonal entries of  $Q$  are non-negative, then we have established that the diagonal entries of  $P$  are upper bounded by 1 (since we know that these entries are non-negative). But the matrix  $Q$  is also a projection matrix. Specifically, we have  $Q^H = (I - P)^H = I - P^H = I - P = Q$  and  $Q^2 = (I - P)(I - P) = I - 2P + P^2 = I - 2P + P = I - P = Q$ . Hence, proceeding exactly as above,

$$Q_{11} = \mathbf{e}_1^H Q \mathbf{e}_1 = \mathbf{e}_1^H Q^H Q \mathbf{e}_1 = \|Q \mathbf{e}_1\|^2. \quad (12)$$

An alternative, less elegant proof starts by observing that  $\mathbf{u}_1, \dots, \mathbf{u}_k$  are an orthonormal basis for a  $k$ -dimensional subspace. Complete this into an orthonormal basis for  $\mathbb{C}^n$  by adding  $\mathbf{u}_{k+1}, \dots, \mathbf{u}_n$ . Then, we can express

$$\mathbf{e}_1 = \sum_{i=1}^n \mu_i \mathbf{u}_i, \quad (13)$$

where  $\mu_i$  are appropriate coefficients satisfying  $\sum_{i=1}^n |\mu_i|^2 = 1$  (since  $\|\mathbf{e}_1\|^2 = 1$ ). Using this representation, we find

$$P_{11} = \|P \mathbf{e}_1\|^2 = \left\| P \sum_{i=1}^n \mu_i \mathbf{u}_i \right\|^2 = \left\| \sum_{i=1}^n \mu_i P \mathbf{u}_i \right\|^2 = \left\| \sum_{i=1}^k \mu_i P \mathbf{u}_i \right\|^2, \quad (14)$$

where the last step is because the vectors  $\mathbf{u}_{k+1}, \dots, \mathbf{u}_n$  are orthogonal to all the vectors  $\mathbf{u}_1, \dots, \mathbf{u}_k$ . Moreover, for any vector  $\mathbf{x}$  inside the subspace spanned by  $\mathbf{u}_1, \dots, \mathbf{u}_k$ , we have  $P\mathbf{x} = \mathbf{x}$ , hence,

$$P_{11} = \left\| \sum_{i=1}^k \mu_i \mathbf{u}_i \right\|^2 = \sum_{i=1}^k |\mu_i|^2 \leq 1, \quad (15)$$

where the second step is because  $\mathbf{u}_1, \dots, \mathbf{u}_k$ , are orthonormal and the last step is because we know that  $\sum_{i=1}^n |\mu_i|^2 = 1$ . This establishes the claim.

**Problem 2: Eckart–Young Theorem**

In class, we proved the converse part of the Eckart–Young theorem for the spectral norm. Here, you do the same for the case of the Frobenius norm.

(a) For any matrix  $A$  of dimension  $m \times n$  and an arbitrary orthonormal basis  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of  $\mathbb{C}^n$ , prove that

$$\|A\|_F^2 = \sum_{k=1}^n \|A\mathbf{x}_k\|^2. \quad (16)$$

(b) Consider any  $m \times n$  matrix  $B$  with  $\text{rank}(B) \leq p$ . Clearly, its null space has dimension no smaller than  $n - p$ . Therefore, we can find an orthonormal set  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n-p}\}$  in the null space of  $B$ . Prove that for such vectors, we have

$$\|A - B\|_F^2 \geq \sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2. \quad (17)$$

(c) (This requires slightly more subtle manipulations.) For any matrix  $A$  of dimension  $m \times n$  and any orthonormal set of  $n - p$  vectors in  $\mathbb{C}^n$ , denoted by  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n-p}\}$ , prove that

$$\sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2 \geq \sum_{j=p+1}^r \sigma_j^2. \quad (18)$$

*Hint:* Consider the case  $m \geq n$  and the set of vectors  $\{\mathbf{z}_1, \dots, \mathbf{z}_{n-p}\}$ , where  $\mathbf{z}_k = V^H \mathbf{x}_k$ . Express your formulas in terms of these and the SVD representation  $A = U\Sigma V^H$ .

(d) Briefly explain how (a)-(c) imply the desired statement.

**Solution 2.** (a) Let us collect the vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  (as columns) into an  $n \times n$  matrix  $X$ . With this, we can express

$$\sum_{k=1}^n \|A\mathbf{x}_k\|^2 = \|AX\|_F^2. \quad (19)$$

Using the result that  $\|A\|_F^2 = \text{trace}(A^H A)$ , we find

$$\|AX\|_F^2 = \text{trace}((AX)^H AX) = \text{trace}(X^H A^H AX) = \text{trace}(A^H AX X^H), \quad (20)$$

where the last step is the property that  $\text{trace}(AB) = \text{trace}(BA)$ . But since by construction,  $X$  is a unitary matrix, we have that  $XX^H$  is simply the identity matrix. Hence,  $\text{trace}(A^H AX X^H) = \text{trace}(A^H A)$ , which completes the proof.

(b) Let us first expand our orthonormal set  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n-p}\}$  to a full basis for  $\mathbb{C}^n$  by including orthonormal vectors  $\{\mathbf{x}_{n-p+1}, \dots, \mathbf{x}_n\}$ . Then, from Part (a), we have

$$\|A - B\|_F^2 = \sum_{k=1}^n \|(A - B)\mathbf{x}_k\|^2 \geq \sum_{k=1}^{n-p} \|(A - B)\mathbf{x}_k\|^2, \quad (21)$$

where the last step is simply because all terms in the sum are non-negative. But by construction,  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n-p}\}$  are in the null space of  $B$ , thus for them,  $B\mathbf{x}_k = \mathbf{0}$ , which implies  $(A - B)\mathbf{x}_k = A\mathbf{x}_k$ . This completes the proof.

(c) The first point of this exercise was to recall the often surprisingly useful “trick” that  $\|\mathbf{y}\|^2 = \text{trace}(\mathbf{y}^H \mathbf{y})$ , where of course the trace-operator does not do anything (yet). Applying this, we can express:

$$\|A\mathbf{x}_k\|^2 = \text{trace}(\mathbf{x}_k^H A^H A \mathbf{x}_k) = \text{trace}(\mathbf{x}_k^H V \Sigma^H U^H U \Sigma V^H \mathbf{x}_k) \quad (22)$$

$$= \text{trace}(\mathbf{z}_k^H \Sigma^H \Sigma \mathbf{z}_k) = \text{trace}(\Sigma^H \Sigma \mathbf{z}_k \mathbf{z}_k^H), \quad (23)$$

where in the last step, we have used the property  $\text{trace}(AB) = \text{trace}(BA)$ . Hence,

$$\sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2 = \sum_{k=1}^{n-p} \text{trace}(\Sigma^H \Sigma \mathbf{z}_k \mathbf{z}_k^H) \quad (24)$$

$$= \text{trace}(\Sigma^H \Sigma \sum_{k=1}^{n-p} \mathbf{z}_k \mathbf{z}_k^H) \quad (25)$$

$$= \sum_{k=1}^n \sigma_k^2 G_{kk}, \quad (26)$$

where the last step is due to the fact that  $\Sigma^H \Sigma$  is a *diagonal* matrix with entries  $\sigma_k^2$ , and where  $G_{ij}$  denote the entries of the matrix  $G = \sum_{k=1}^{n-p} \mathbf{z}_k \mathbf{z}_k^H$ . The matrix  $G$  is a projection matrix. As we have seen in an earlier homework problem, its trace is  $n - p$  and its diagonal entries are non-negative and no larger than one. Under these constraints, it should be clear that the last expression is minimized if we select  $G_{11} = G_{22} = \dots = G_{pp} = 0$  and  $G_{p+1,p+1} = \dots = G_{nn} = 1$ . Hence,

$$\sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2 = \sum_{k=1}^n \sigma_k^2 G_{kk} \quad (27)$$

$$\geq \sum_{k=p+1}^n \sigma_k^2 \quad (28)$$

$$\geq \sum_{k=p+1}^r \sigma_k^2, \quad (29)$$

which completes the proof.

(d) Combining Parts (b) and (c):

$$\|A - B\|_F^2 \geq \sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2 \geq \sum_{j=p+1}^r \sigma_j^2 \quad (30)$$

shows that for *any* matrix  $B$  of rank  $p$ , we have the above lower bound. This is precisely the statement needed to complete the proof of the Eckart-Young theorem for the Frobenius norm.

*Additional remark:* Another proof of the Eckart-Young theorem (which works both for the Frobenius and the spectral norm) leverages the Weyl theorem, which states that for any two matrices  $C$  and  $D$  of the same dimension ( $m \times n$ , and assume w.l.o.g.  $m \geq n$ ), we have that

$$\sigma_{i+j-1}(C + D) \leq \sigma_i(C) + \sigma_j(D), \quad \text{for } 1 \leq i, j \leq n, \text{ and } i + j - 1 \leq n. \quad (31)$$

I am not aware of a simple proof of this theorem (the standard proof uses the variational characterization of eigenvalues). But suppose that  $B$  is of rank no larger than  $k$ , meaning that  $\sigma_i(B) = 0$  for  $i > k$ . Then, setting  $C = A - B$  and  $D = B$ , Weyl's theorem says that

$$\sigma_{i+k}(A) \leq \sigma_i(A - B) \quad \text{for } 1 \leq i \leq n - k, \quad (32)$$

and thus,

$$\|A - B\|_F^2 \geq \sum_{i=1}^{n-k} \sigma_i^2(A - B) \geq \sum_{i=k+1}^n \sigma_i^2(A). \quad (33)$$

### Problem 3: A Hilbert space of matrices

In this problem, we consider the set of matrices  $A \in \mathbb{R}^{m \times n}$  with standard matrix addition and multiplication by scalar.

(a) Briefly argue that this is indeed a vector space, using the definition given in class.

(b) Show that  $\langle A, B \rangle = \text{trace}(B^H A)$  is a valid inner product.

(c) Explicitly state the norm induced by this inner product. Is this a norm that you have encountered before?

(d) Consider as a further inner product candidate the form  $\langle A, B \rangle = \text{trace}(B^H W A)$ , where  $W$  is a square ( $m \times m$ ) matrix. Give conditions on  $W$  such that this is a valid inner product. Explicit and detailed arguments are required for full credit.

**Solution 3.** *Note: In the following, we solve assuming the more general case of complex valued matrices. It should be easier to solve for the real-valued case.*

(a) We need to check some properties. Because the space is that of matrices,

1. Commutativity holds.
2. Associativity holds.
3. Distributivity holds.
4. The 0 element is the all 0's matrix  $\mathbf{0} \in \mathbb{C}^{m \times n}$ .
5. For all  $A \in \mathbb{C}^{m \times n}$ , we have that the element  $-A \in \mathbb{C}^{m \times n}$  is such that  $A + (-A) = \mathbf{0}$ .
6. For all  $A \in \mathbb{C}^{m \times n}$ , we have that  $I_{m \times m} A = A$ .

So this is indeed a vector space.

(b) Here, we check the properties of an inner product space. Letting  $A, B, C \in \mathbb{C}^{m \times n}, \alpha \in \mathbb{C}$ , we have that

1.  $\langle A + C, B \rangle = \text{Tr}(B^H(A + C)) = \text{Tr}(B^H A + B^H C) = \text{Tr}(B^H A) + \text{Tr}(B^H C) = \langle A, B \rangle + \langle C, B \rangle$ , where we used the linearity of the trace operator.
2.  $\langle \alpha A, B \rangle = \text{Tr}(B^H \alpha A) = \alpha \text{Tr}(B^H A) = \alpha \langle A, B \rangle$ , where we used the linearity of the trace operator.
3.  $\langle A, B \rangle = \text{Tr}(B^H A) = \text{Tr}((A^H B)^H) = \text{Tr}(A^H B)^* = \langle B, A \rangle^*$ , where we used the linearity of the trace operator and that conjugation is also linear.

4. We want  $\langle A, A \rangle = \text{Tr}(A^H A) \geq 0$ . Since  $A^H A$  is normal, it is also positive semi-definite, and so all its eigenvalues are positive. One of the property of the trace is that it is equal to the sum of eigenvalues of the matrix considered. In our case, this means  $\text{Tr}(A^H A) \geq 0$  since the eigenvalues are all non-negative.

(c) We have that the norm is  $\sqrt{\langle A, A \rangle} = \sqrt{\text{Tr}(A^H A)} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2}$ , which we recognize to be the Frobenius norm. So  $\sqrt{\langle A, A \rangle} = \|A\|_F$ .

(d) We check that the properties of an inner product space hold, and add conditions on  $W$  when necessary.

1.  $\langle A+C, B \rangle = \text{Tr}(B^H W(A+C)) = \text{Tr}(B^H W A + B^H W C) = \text{Tr}(B^H A) + \text{Tr}(B^H C) = \langle A, B \rangle + \langle C, B \rangle$ , so no restriction on  $W$  necessary here.
2.  $\langle \alpha A, B \rangle = \text{Tr}(B^H W \alpha A) = \alpha \text{Tr}(B^H W A) = \alpha \langle A, B \rangle$ , so no restriction on  $W$  necessary here.
3. On one side we have  $\langle A, B \rangle = \text{Tr}(B^H W A)$ . On the other side we have  $\langle B, A \rangle^* = \text{Tr}(A^H W B)^* = \text{Tr}((A^H W B)^H) = \text{Tr}(B^H W^H A)$ . To have both sides equal, we need  $W = W^H$ , i.e.,  $W$  should be Hermitian.
4. We want  $\langle A, A \rangle = \text{Tr}(A^H W A) \geq 0$ . That is we would like  $A^H W A$  to be positive semi-definite. By definition, this would mean that for any  $\mathbf{z} \in \mathbb{C}^n$ , we want  $\mathbf{z}^H A^H W A \mathbf{z} \geq 0$ . Now note that  $A \mathbf{z}$  is just another vector, so we can write that we want  $\mathbf{z}^H A^H W A \mathbf{z} = (A \mathbf{z})^H W (A \mathbf{z}) \geq 0$  for all  $\mathbf{z}$ . So we conclude that this is the same as asking that  $W$  is positive semi-definite.

Hence, for the inner product  $\langle A, B \rangle = \text{Tr}(B^H W A)$  to be valid, we need  $W$  to be Hermitian and positive semi-definite.

#### Problem 4: Haar Wavelet

*This problem is taken from Vetterli/Kovacevic, p. 295.*

Consider the wavelet series expansion of continuous-time signals  $f(t)$  and assume that  $\psi(t)$  is the Haar wavelet.

(a) Give the expansion coefficients for  $f(t) = 1, t \in [0, 1]$ , and 0 otherwise.

(b) Verify that for  $f(t)$  as in Part (a),  $\sum_m \sum_n \|\langle \psi_{m,n}, f \rangle\|^2 = 1$  (i.e., Parseval's identity).

**Solution 4.** a) We have that  $a_{m,n} = \int_{-\infty}^{\infty} f(t') \psi_{m,n}(t') dt' = 2^{-m/2} \int_0^1 \psi(2^{-m}t' - n) dt' = 2^{m/2} \int_{-n}^{-n+2^{-m}} \psi(\tau) d\tau$  where we used the substitution  $\tau := 2^{-m}t' - n$ . We can see that the coefficients can be computed by simply integrating over the mother wavelet  $\psi$ .

First assume that  $n = 0$ . Then, for  $m = 0$ ,  $a_{0,0} = 0$  due to symmetric integration of the Haar wavelet around its midpoint. For  $m \geq 1$ , it can be easily verified that  $a_{m,0} = 2^{-m/2}$ . For  $m \leq -1$ , the coefficients are zero because we integrate over the whole support of the Haar wavelet.

Now assume that  $n \leq -1$ . Then we have that  $-n \geq 1$ , hence our integration interval misses the support of the mother wavelet towards the right and therefore  $a_{m,n} = 0$ .

Likewise, for  $n \geq 1$  we always have that either  $2^{-m} - n \leq 0$  (in which case we miss the support towards the left) or  $2^{-m} + n \geq 1$  (we integrate over the whole wavelet) and hence  $a_{m,n} = 0$ .

In conclusion, we have that  $a_{m,n} = 2^{-m/2}$  for  $(n = 0) \wedge (m \geq 1)$  and  $a_{m,n} = 0$  everywhere else.

b)  $\sum_{m,n} |a_{m,n}|^2 = \sum_{m=1}^{\infty} |a_{m,0}|^2 = \sum_{m=0}^{\infty} (\frac{1}{2})^m - 1 = \frac{1}{1-\frac{1}{2}} - 1 = 1$

**Problem 5: Dual Representation of Norm**

(a) Assume that  $p > 0$  and  $q > 0$  fulfils  $1/p + 1/q = 1$ - Show that the following inequality holds for all  $a \geq 0$  and  $b \geq 0$ .

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q} \tag{34}$$

Show that the equality holds if  $a^p = b^q$ . [Hint: Use the concavity of log function]

(b) Given vectors  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^n$  show that,

$$\frac{\sum_{i=1}^n |x_i y_i|}{\|\mathbf{x}\|_p \|\mathbf{y}\|_q} \leq 1 \tag{35}$$

What is the condition for equality?

(c) Show that

$$\|\mathbf{x}\|_p = \sup_{\|\mathbf{y}\|_{p^*} = 1} \langle \mathbf{y}, \mathbf{x} \rangle; \mathbf{y} \in \mathbb{R}^n, \|\mathbf{y}\|_{p^*} = 1.$$

where  $1/p + 1/p^* = 1$

**Solution 5.** (a) Taking the log of the right hand side we have

$$\begin{aligned} \log\left(\frac{a^p}{p} + \frac{b^q}{q}\right) &\geq \frac{1}{p} \log(a^p) + \frac{1}{q} \log(b^q) && \text{(Concavity of log, } \frac{1}{p} + \frac{1}{q} = 1) \\ &= \log(a) + \log(b) = \log(ab) \end{aligned}$$

If  $a^p = b^q = c$ , the inequality becomes equality as we have  $\log\left(\left(\frac{1}{p} + \frac{1}{q}\right)c\right) \geq \left(\frac{1}{p} + \frac{1}{q}\right) \log(c)$

(b) Define  $\tilde{x}_i = \frac{|x_i|}{\|\mathbf{x}\|_p}$  and  $\tilde{y}_i = \frac{|y_i|}{\|\mathbf{y}\|_q}$  for  $i \in [1, n]$ . Note that

$$\sum_{i=1}^n \tilde{x}_i \tilde{y}_i = \frac{\sum_{i=1}^n |x_i y_i|}{\|\mathbf{x}\|_p \|\mathbf{y}\|_q}$$

From the result of part (a), we have

$$\begin{aligned} \sum_{i=1}^n \tilde{x}_i \tilde{y}_i &\leq \frac{1}{p} \sum_{i=1}^n \tilde{x}_i^p + \frac{1}{q} \sum_{i=1}^n \tilde{y}_i^q \\ &= \frac{1}{p} \frac{\sum_{i=1}^n |x_i|^p}{\|\mathbf{x}\|_p^p} + \frac{1}{q} \frac{\sum_{i=1}^n |y_i|^q}{\|\mathbf{y}\|_q^q} && \left(\sum_{i=1}^n |x_i|^p = \|\mathbf{x}\|_p^p\right) \\ &= \frac{1}{p} + \frac{1}{q} = 1 \end{aligned}$$

As shown in part (a), equality happens only if  $\tilde{x}_i^p = \tilde{y}_i^q \forall i \in [1, n]$ .

(c)

By using the inequality proved in part (b) we have that for every  $\mathbf{y}$  such that  $\|\mathbf{y}\|_{p^*} = 1$ ,

$$\langle \mathbf{y}, \mathbf{x} \rangle \leq \sum_{i=1}^n |x_i y_i| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_{p^*} = \|\mathbf{x}\|_p.$$

In order to achieve equality, we know from part (b) that we must have  $\tilde{x}_i^p = \tilde{y}_i^q$  or  $(\frac{|x_i|}{\|\mathbf{x}\|_p})^p = (\frac{|y_i|}{\|\mathbf{y}\|_{p^*}})^{p^*} \forall i \in [1, n]$ .

Since,  $\|\mathbf{y}\|_{p^*} = 1$ , we choose  $y_i$  as

$$y_i = \text{sign}(x_i) \frac{|x_i|^{p/p^*}}{(\sum_{i=1}^n |x_i|^p)^{\frac{1}{p^*}}}$$

### Problem 6: Finding the Fair Coin

#### To be done in the exercise session on Oct 17.

Your colleague challenges you to a game. It goes as follows: You are given three coins, one being fair, two being weighted. Coin one has a probability of flipping heads  $1/2$ , coin two has probability of flipping heads  $1/2 - p$ , and coin three has probability of flipping heads  $1/2 + p$ . You can assume that  $p \in (0, 1/3]$ . You are allowed to flip every coin  $m$  times. If thereafter, you manage to correctly identify the fair coin, you win. Otherwise you lose.

- Describe a simple strategy that is winning for  $m \rightarrow \infty$ .
- Write down the events that make this strategy fail (assume finite  $m$ ).
- Give a simple upper bound on the total failure probability in terms of  $m$  and  $p$ .
- Your colleague makes the following proposal: you can flip each coin  $m = 20$  times, and he ensures you that  $p = 1/3$  (he is a good friend, so you trust him). If you can identify the fair coin, you win 3 CHF, otherwise you lose 2 CHF. Do you accept?

In general, for  $p \in (0, 1/3]$  fixed, for which values of  $m$ ,  $\alpha$  would you agree to play the game, where  $\alpha$  is defined as the ratio between the potential gain and potential loss in CHF? You can assume that you know  $p$ . It is sufficient to state bounds that are tight up to multiplicative constants.

**Solution 6.** a) Associate heads with 0s and tails with 1s. Compute the empirical means  $\hat{\mu}_1, \hat{\mu}_2, \hat{\mu}_3$  of each of the three coins. Select the coin corresponding to the second largest empirical mean. This strategy is winning for  $m \rightarrow \infty$  by the law of large numbers.

$$b) E_1 = \{\hat{\mu}_2 < \hat{\mu}_3 < \hat{\mu}_1\}, E_2 = \{\hat{\mu}_1 < \hat{\mu}_2 < \hat{\mu}_3\}, E_3 = \{\hat{\mu}_3 < \hat{\mu}_2 < \hat{\mu}_1\}, E_4 = \{\hat{\mu}_1 < \hat{\mu}_3 < \hat{\mu}_2\}$$

c) Denote by  $X_{j,k}$  the random variable associated with flipping the  $k$ -th coin for the  $j$ -th time. Then,

$$\mathbb{P}(\{\text{failure}\}) = \mathbb{P}\left(\bigcup_i E_i\right) \tag{36}$$

$$\leq \sum_i \mathbb{P}(E_i) \tag{37}$$

$$\leq 4 \cdot \mathbb{P}(\hat{\mu}_3 < \hat{\mu}_1) \tag{38}$$

$$= 4 \cdot \mathbb{P}\left(\frac{1}{m} \sum_j X_{j,3} - X_{j,1} \leq 0\right) \tag{39}$$

$$= 4 \cdot \mathbb{P}\left(\frac{1}{m} \sum_j X_{j,3} - (1/2 + p) - X_{j,1} + 1/2 \leq -p\right) \tag{40}$$

$$\leq 4e^{-mp^2/4} \tag{41}$$



Where in the last inequality we used Hoeffding's bound for centered, bounded random variables.

*d)* Plugging in  $m = 20$ ,  $p = 1/3$  into the above bound we obtain that  $\mathbb{P}(\{\text{failure}\}) \leq 2.29$  which is vacuous. Hence we have no guarantee that will make a gain with the above mentioned strategy and we should not accept. In general, we should accept whenever  $\mathbb{P}(\{\text{failure}\})/\mathbb{P}(\{\text{success}\}) < \alpha$ , i.e. when  $4e^{-mp^2/4} < 1/(1 - 4^{-mp^2/4}) < \alpha \Leftrightarrow m > 4 \ln(4/\alpha + 4)/p^2$ . Note that sharper results can be obtained by bounding  $\mathbb{P}(\{\text{failure}\})$  more carefully.