

Theoretical task 1

All solutions should be short, mathematically precise and contain proof unless qualitative explanation/intuition is needed.

- Suppose $x \in \mathbb{R}^D$ is a feature vector. Prove that *whitening transformation* $f = \Sigma^{-1/2}(x - \mu)$, where $\mu = \mathbb{E}x$, $\Sigma = \text{cov}[x, x]$, will give new feature vector f with properties:

- $\mathbb{E}f = \mathbf{0}$ (all zeroes vector)
- $\text{cov}[f, f] = I$ (identity matrix)

- Consider training set x_1, \dots, x_N and some linear subspace L_K with lower dimensionality $K \leq D$. Let $x_i = p_i + h_i$ where p_i are projections of x_i onto L_K and h_i are orthogonal complements. Suppose we perform optimization over different K -dimensional subspaces L_K . Prove equivalence of the following two optimization tasks:

- $\sum_{i=1}^N \|h_i\|^2 \rightarrow \min_{L_K}$
- $\sum_{i=1}^N \|p_i\|^2 \rightarrow \max_{L_K}$

Comment: $\|z\|$ is L_2 norm of vector z .

- Write stochastic gradient descent with minibatch size=1 for the following losses:

- $\mathcal{L}(M) = e^{-M}$
- $\mathcal{L}(M) = [1 - M]_+$

Why classification quality (evaluated by by margin) on object from the minibatch cannot decrease?

- Prove that $K(x, x') = e^{-\gamma \langle x - x', x - x' \rangle}$, $\gamma > 0$ is a Mercer kernel.

Hint: use operations generating new kernels out of existing kernels.

- Consider a binary classifier $\hat{y}(x) = \text{sign}(g(x) - \mu)$ with discriminant function $g(x)$ and some threshold μ . Suppose you know $TPR(\mu)$ and $FPR(\mu)$. Now consider an inverted classifier $\tilde{y}(x) = \text{sign}(\mu - g(x))$. Write out $TPR(\mu)$ and $FPR(\mu)$ measures for it in terms of original classifier. Explain.
- Consider multiclass classification performed by M classifiers $f_1(x), \dots, f_M(x)$. Let probability of mistake be constant $p \in (0, \frac{1}{2})$: $p(f_m(x) \neq y) = p \forall m$ and suppose all models make mistakes or correct guesses independently of each other. Let $F(x)$ be majority voting aggregation function (voting for most popular class among predicted by $f_1(x), \dots, f_M(x)$). Prove that $\forall(x, y) p(F(x) \neq y) \rightarrow 0$ as $M \rightarrow \infty$.

Hint: use central limit theorem, consider fraction of errors. What can be said about its expectation and variance?