

1 Theoretical Part

Based on Lecture 6 and Recitation 10

1. Let $k(\mathbf{x}, \mathbf{x}')$ be a valid PSD kernel. Provide a valid PSD kernel $\tilde{k}(\mathbf{x}, \mathbf{x}')$, constructed from k , which is guaranteed to be normalized. That is, for all \mathbf{x} it holds that $\tilde{k}(\mathbf{x}, \mathbf{x}) = 1$. Prove your answer.

Let $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as

$$\tilde{k}(x, x') = k(x, x') / \sqrt{k(x, x) \cdot k(x', x')}$$

First, we want to demonstrate that $\tilde{k}(x, x')$ is a valid positive semidefinite kernel. We can observe that $\tilde{k}(x, x')$ is symmetric, we need to show that it is PSD and not negative and normalized.

According to Mercer's theorem, because $k(x, x')$ is a valid kernel, we can express it in a specific form:

$$K(x, x') = (\Phi(x)(\text{transpose})) \Phi(x')$$

Therefore we can conclude:

$$\tilde{k}(x, x') = k(x, x') / \sqrt{k(x, x) \cdot k(x', x')} =$$

$$((\Phi(x)(\text{transpose})) \Phi(x')) / (\sqrt{(\Phi(x)(\text{transpose})) \Phi(x)} \cdot \sqrt{(\Phi(x')(\text{transpose})) \Phi(x')}) =$$

$$((\Phi(x)(\text{transpose})) \Phi(x')) / (\text{norma}(\Phi(x)) \cdot \text{norma}(\Phi(x'))) = ((\tilde{\Phi}(x)) \text{transpose}) * (\tilde{\Phi}(x'))$$

Where $\tilde{\Phi}(x) = \Phi(x) / \text{norma}(\Phi(x))$.

And now we showed that \tilde{k} is normalized and also PSD.

2. Consider a data set $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{\pm 1\}$, and a feature map $\psi: \mathbb{R}^d \rightarrow \mathcal{F}$ where \mathcal{F} is some feature space. Give an example of a data set S and a feature map ψ such that S is not linearly separable in \mathbb{R}^d (for $d \geq 2$) but that the transformed data set $S_\psi = \{(\psi(\mathbf{x}_i), y_i)\}_{i=1}^m$ is linearly separable in \mathcal{F} .

Given the dataset observed during the recitation, specifically involving two co-centric rings with varying radii, by the following mapping function we can transform and achieved into a linearly separable dataset in \mathbb{R}^d : $\psi(x_1, x_2) = (x_1, x_2, x_1^2 + x_2^2)$

[Type here]

3. $k_1(\mathbf{x}, \mathbf{y})$ and $k_2(\mathbf{x}, \mathbf{y})$ are valid kernels, then:

$$k_{\times}(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y})$$

is also a valid kernel. To prove this we'll use the fact that valid kernels are positive semi-definite.

You may find the following identities helpful (but don't have to use them):

$$\mathbf{x}^\top \mathbf{A} \mathbf{y} = \text{Tr}[\mathbf{x}^\top \mathbf{A} \mathbf{y}] = \text{Tr}[\mathbf{y} \mathbf{x}^\top \mathbf{A}] \quad (1)$$

$$\text{Tr}[\mathbf{A} \mathbf{B}] = \sum_i [\mathbf{A} \mathbf{B}]_{ii} = \sum_i \sum_j A_{ij} B_{ji} \quad (2)$$

where \mathbf{x} and \mathbf{y} are vectors while \mathbf{A} and \mathbf{B} are matrices.

To prove that the kernel $k_{\times}(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y})$ is also a valid kernel, we need to show that it satisfies the positive semidefinite property.

Let \mathbf{K}_1 and \mathbf{K}_2 be the kernel matrices corresponding to $k_1(\mathbf{x}, \mathbf{y})$ and $k_2(\mathbf{x}, \mathbf{y})$, respectively. The kernel matrix \mathbf{K} for $k_{\times}(\mathbf{x}, \mathbf{y})$ is given by $\mathbf{K} = \mathbf{K}_1 \cdot \mathbf{K}_2$.

To demonstrate that \mathbf{K} is a valid kernel, we must show that for any vector \mathbf{x} , the quadratic form $\mathbf{x}' \mathbf{K} \mathbf{x}$ is non-negative.

Considering $\mathbf{x}' \mathbf{K} \mathbf{x}$: $\mathbf{x}' \mathbf{K} \mathbf{x} = \mathbf{x}' (\mathbf{K}_1 \cdot \mathbf{K}_2) \mathbf{x}$

Using the provided identities, we can simplify this expression:

$$\mathbf{x}' \mathbf{K} \mathbf{x} = \text{Trace}[\mathbf{x}' (\mathbf{K}_1 \cdot \mathbf{K}_2) \mathbf{x}] = \sum [\mathbf{x}' (\mathbf{K}_1 \cdot \mathbf{K}_2) \mathbf{x}]_{ii}$$

$$\text{Expanding } \mathbf{K}_1 \cdot \mathbf{K}_2: \mathbf{x}' \mathbf{K} \mathbf{x} = \sum [\mathbf{x}' \mathbf{K}_1 \mathbf{K}_2 \mathbf{x}]_{ii}$$

Since \mathbf{K}_1 and \mathbf{K}_2 are valid kernel matrices, they can be expressed as $\mathbf{K}_1 = \mathbf{A}_1' \mathbf{A}_1$ and $\mathbf{K}_2 = \mathbf{A}_2' \mathbf{A}_2$, where \mathbf{A}_1 and \mathbf{A}_2 are matrices.

$$\text{Substituting these expressions, we get: } \mathbf{x}' \mathbf{K} \mathbf{x} = \sum [\mathbf{x}' (\mathbf{A}_1' \mathbf{A}_1 \cdot \mathbf{A}_2' \mathbf{A}_2) \mathbf{x}]_{ii}$$

Let's define a new matrix $\mathbf{A} = \mathbf{A}_1 \mathbf{A}_2$. Since \mathbf{A}_1 and \mathbf{A}_2 are positive semidefinite, $\mathbf{A} = \mathbf{A}_1 \mathbf{A}_2$ is also positive semidefinite.

$$\text{We can rewrite the expression as: } \mathbf{x}' \mathbf{K} \mathbf{x} = \sum [\mathbf{x}' \mathbf{A}' \mathbf{A} \mathbf{x}]_{ii}$$

Since \mathbf{A} is positive semidefinite, the quadratic form $\mathbf{x}' \mathbf{A}' \mathbf{A} \mathbf{x}$ is non-negative for any vector \mathbf{x} .

Therefore, we have shown that $\mathbf{x}' \mathbf{K} \mathbf{x}$ is non-negative, proving that the kernel $k_{\times}(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y})$ is a valid kernel.

3. $k_1(\mathbf{x}, \mathbf{y})$ and $k_2(\mathbf{x}, \mathbf{y})$ are valid kernels, then:

$$k_{\times}(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y})$$

is also a valid kernel. To prove this we'll use the fact that valid kernels are positive semi-definite.

You may find the following identities helpful (but don't have to use them):

$$\mathbf{x}^T \mathbf{A} \mathbf{y} = \text{Tr}[\mathbf{x}^T \mathbf{A} \mathbf{y}] = \text{Tr}[\mathbf{y} \mathbf{x}^T \mathbf{A}] \quad (1)$$

$$\text{Tr}[\mathbf{A} \mathbf{B}] = \sum_i [\mathbf{A} \mathbf{B}]_{ii} = \sum_i \sum_j \mathbf{A}_{ij} \mathbf{B}_{ji} \quad (2)$$

where \mathbf{x} and \mathbf{y} are vectors while \mathbf{A} and \mathbf{B} are matrices.

- (a) Let $k(\mathbf{x}, \mathbf{y})$ be a valid kernel and suppose that \mathbf{K} is the kernel's Gram matrix over some finite set of points $\{\mathbf{x}_i\}_{i=1}^N$, such that $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. Show that for any finite set of N points, there exists some function $f: \mathcal{X} \mapsto \mathbb{R}^N$ such that:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{f}^T(\mathbf{x}_i) \mathbf{f}(\mathbf{x}_j) \quad (3)$$

where \mathcal{X} is space of the points \mathbf{x}_i . Using this fact, show that:

$$k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y}) = \sum_i \sum_j g_i(\mathbf{x}) f_j(\mathbf{x}) f_j(\mathbf{y}) g_i(\mathbf{y}) \quad (4)$$

where $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are valid kernels, and some functions $f, g: \mathcal{X} \mapsto \mathbb{R}^N$, where $f_i(\mathbf{x})$ denotes the i^{th} index of the output of $f(\mathbf{x})$.

- (b) Conclude that $k_{\times}(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y}) = \mathbf{h}^T(\mathbf{x}) \cdot \mathbf{h}(\mathbf{y})$ for some function $\mathbf{h}(\cdot)$, thereby proving that $k_{\times}(\cdot, \cdot)$ is a valid kernel.

3.3.a

Given a valid kernel $k(\mathbf{x}, \mathbf{y})$ and its Gram matrix \mathbf{K} computed from a finite set of N points $\{\mathbf{x}_i\}$, we want to show that there exists a function $\mathbf{f}: \mathcal{X} \rightarrow \mathbb{R}^N$ such that $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{f}^T(\mathbf{x}_i))\mathbf{f}(\mathbf{x}_j)$, where \mathcal{X} represents the space of the points \mathbf{x}_i .

By using the eigendecomposition of the Gram matrix \mathbf{K} as $\mathbf{K} = \mathbf{V} \mathbf{S} \mathbf{V}^T$, where \mathbf{V} is the matrix of eigenvectors and \mathbf{S} is the diagonal matrix of eigenvalues, we can define $\mathbf{f}(\mathbf{x}) = [\sqrt{\lambda_1} \mathbf{v}_1(\mathbf{x}), \sqrt{\lambda_2} \mathbf{v}_2(\mathbf{x}), \dots, \sqrt{\lambda_N} \mathbf{v}_N(\mathbf{x})]$, where $\mathbf{v}_i(\mathbf{x})$ represents the i -th eigenvector associated with the i -th eigenvalue λ_i .

Calculating $(\mathbf{f}^T(\mathbf{x}_i))\mathbf{f}(\mathbf{x}_j)$ simplifies to $\sqrt{\lambda_1} \mathbf{v}_1(\mathbf{x}_i)^T \mathbf{v}_1(\mathbf{x}_j) + \sqrt{\lambda_2} \mathbf{v}_2(\mathbf{x}_i)^T \mathbf{v}_2(\mathbf{x}_j) + \dots + \sqrt{\lambda_N} \mathbf{v}_N(\mathbf{x}_i)^T \mathbf{v}_N(\mathbf{x}_j)$. Since the eigenvectors are orthogonal, $\mathbf{v}_i^T(\mathbf{x}_i) \mathbf{v}_j(\mathbf{x}_j)$ equals 1 if $i = j$, and 0 if $i \neq j$.

Thus, we can simplify $(\mathbf{f}^T(\mathbf{x}_i))\mathbf{f}(\mathbf{x}_j)$ to $\sqrt{\lambda_1} K_{1j} + \sqrt{\lambda_2} K_{2j} + \dots + \sqrt{\lambda_N} K_{Nj}$, where K_{ij} represents the entries of the Gram matrix \mathbf{K} .

Therefore, by defining $\mathbf{f}(\mathbf{x})$ as mentioned, we have shown that for any finite set of N points, there exists a function $\mathbf{f}: \mathcal{X} \rightarrow \mathbb{R}^N$ such that $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{f}^T(\mathbf{x}_i))\mathbf{f}(\mathbf{x}_j)$, where \mathcal{X} represents the space of the points \mathbf{x}_i .

3.4.a

To show that $k_1(x, y) \cdot k_2(x, y)$ can be written as $\sum_i \sum_j g_i(x) f_j(x) f_j(y) g_i(y)$, where $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are valid kernels and $f, g: X \rightarrow \mathbb{R}^N$ are functions, we will use the result that states for any valid kernel $k(x, y)$, there exist functions $f: X \rightarrow \mathbb{R}^N$ and $g: X \rightarrow \mathbb{R}^N$ such that $k(x, y) = (f^T(x))g(y)$.

Let's consider the expression $k_1(x, y) \cdot k_2(x, y)$. Using the functions f and g , we can rewrite it as $(f^T(x))g(y)$ multiplied by itself.

Expanding this expression, we have $(f^T(x))^2 g(y)^2$. Since $(f^T(x))^2$ is a scalar value, we can bring it out of the summation.

Next, we can express $g(y)^2$ as a matrix $G = g(y) g(y)^T$, where G is an $N \times N$ symmetric matrix.

Now, we can rewrite the expression as $(f^T(x))^2 G$.

Expanding the expression $(f^T(x))^2 G$, we obtain $\sum_i \sum_j g_i(x) f_j(x) f_j(y) g_i(y)$, where $g_i(x)$ represents the i -th component of the output of $g(x)$, and $f_j(x)$ represents the j -th component of the output of $f(x)$.

Therefore, we have shown that $k_1(x, y) \cdot k_2(x, y)$ can be represented as $\sum_i \sum_j g_i(x) f_j(x) f_j(y) g_i(y)$.

.b

By our previous Q, we have shown that $k_1(x, y) \cdot k_2(x, y)$ can be expressed as:

$$k_1(x, y) \cdot k_2(x, y) = \sum_i \sum_j g_i(x) \cdot f_j(x) \cdot f_j(y) \cdot g_i(y),$$

where i represents the index running on the first sigma and j represents the index running on the second sigma.

Now, let's define a new function $h(x)$ as: $h(x) = [g_1(x) \cdot f_1(x), g_2(x) \cdot f_2(x), \dots, g_i(x) \cdot f_j(x), \dots]$.

Considering the inner product of the transpose of $h(x)$ and $h(y)$, we find:

$$(h(x))^T \cdot h(y) = [g_1(x) \cdot f_1(x), g_2(x) \cdot f_2(x), \dots, g_i(x) \cdot f_j(x), \dots]^T \cdot [g_1(y) \cdot f_1(y), g_2(y) \cdot f_2(y), \dots, g_i(y) \cdot f_j(y), \dots]$$

Expanding the inner product, we have: $(h(x))^T \cdot h(y) = \sum_i \sum_j (g_i(x) \cdot f_j(x)) \cdot (g_i(y) \cdot f_j(y))$ which is equivalent to the expression we obtained earlier for $k_1(x, y) \cdot k_2(x, y)$.

Therefore, we have shown that $k_1(x, y) \cdot k_2(x, y)$ can be expressed as the inner product of the transpose of $h(x)$ and $h(y)$, where $h(x)$ is defined as $h(x) = [g_1(x) \cdot f_1(x), g_2(x) \cdot f_2(x), \dots, g_i(x) \cdot f_j(x), \dots]$.

[Type here]

1.2 PCA

Based on Lecture 9 and Recitation 11

- Let $X : \Omega \rightarrow \mathbb{R}^d$ be a random variable with zero mean and covariance $\Sigma \in \mathbb{R}^{d \times d}$. Show that for any $\mathbf{v} \in \mathbb{R}^d$, where $\|\mathbf{v}\|_2 = 1$, the variance of $\langle \mathbf{v}, X \rangle$ is not larger than variance obtained by the PCA embedding of X into a one-dimension subspace (assume that the PCA uses the actual Σ).

Let $\mathbf{v} \in \mathbb{R}^d$ such that $\|\mathbf{v}\|^2 = 1$ and denote $X^\sim \equiv \langle \mathbf{v}, X \rangle$. We want to calculate the expectation and variance of X^\sim .

Expectation: $E[X^\sim] = \langle \mathbf{v}, E[X] \rangle = 0$, since the projection of $E[X]$ onto \mathbf{v} is zero.

Variance: $\text{Var}[X^\sim] = E[X^{\sim 2}] = E[(\langle \mathbf{v}, X \rangle)^2] = E[\langle \mathbf{v}, X \rangle \cdot \langle \mathbf{v}, X \rangle] = \langle \mathbf{v}, X \rangle \cdot \langle \mathbf{v}, X \rangle = (\mathbf{v}^T \Sigma \mathbf{v})$ (using covariance matrix Σ) $\leq (\mathbf{u}_1^T \Sigma \mathbf{u}_1)$ (as \mathbf{v} is a unit vector, and \mathbf{u}_1 is the leading eigenvector of Σ) $= \lambda_1$ (since λ_1 is the largest eigenvalue of Σ)

At the end, the expectation of X^\sim is $E[X^\sim] = 0$, and the variance of X^\sim is $\text{Var}[X^\sim] \leq \lambda_1$, where λ_1 is the largest eigenvalue of the covariance matrix Σ .

1.3 Convex optimization

Based on Lecture 11 and Recitations 2,12

- Let $f_1, \dots, f_m : C \rightarrow \mathbb{R}$ be a set of convex functions and $\gamma_1, \dots, \gamma_m \in \mathbb{R}_+$. Prove from definition that $g(\mathbf{u}) = \sum_{i=1}^m \gamma_i f_i(\mathbf{u})$ is a convex function.
- Give a counterexample for the following claim: Given two functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$, define a new function $h : \mathbb{R} \rightarrow \mathbb{R}$ by $h = f \circ g$. If f and g are convex then h is convex as well.
- Let $f : C \rightarrow \mathbb{R}$ be a function defined over a convex set C . Prove that f is convex iff its *epigraph* is a convex set, where $\text{epi}(f) = \{(u, t) : f(u) \leq t\}$.
- Let $f_i : V \rightarrow \mathbb{R}, i \in I$. Let $f : V \rightarrow \mathbb{R}$ given by

$$f(u) = \sup_{i \in I} f_i(u).$$

If f_i are convex for every $i \in I$, then f is also convex.

1.3.1

We know that for each i belong to $[m]$ f_i is a convex function.

By definition we get a belong $[0,1]$

$aF(u) + (1-a)F(v) \geq F(au + (1-a)v)$ to each u,v belongs to C

multiply positive constant q we get:

$a \cdot q \cdot F(u) + (1-a) \cdot q \cdot F(v) \geq q \cdot F(au + (1-a)v)$

so we conclude that for each i $G_i = q_i f_i(u)$ is a convex function

now we want to conclude it on the sum

we notice that: (the index on each sigma is i from 1 to m) $\sum a G_i(u) + (1-a) G_i(v) = \sum [a G_i(u) + (1-a) G_i(v)] \geq \sum G_i(au + (1-a)v)$.

1.3.2

Counterexample: Let's consider a counterexample to the claim that if f and g are convex functions, then $h = f \circ g$ is also convex.

Take $f(x) = x^2$ and $g(x) = |x|$. Both f and g are convex functions individually, but let's examine the composition $h = f \circ g$.

$$h(x) = f(g(x)) = f(|x|) = |x|^2 = x^2.$$

The function $h(x) = x^2$ is not convex. To see this, we can examine the second derivative of $h(x)$: $h''(x) = 2$,

which is a constant. Since the second derivative is positive (non-negative) everywhere, $h(x)$ does not satisfy the definition of convexity.

Therefore, we have provided a counterexample where f and g are convex functions, but the composition $h = f \circ g$ is not convex.

1.3.3

Assume f is convex. We need to prove that its epigraph, $\text{epi}(f)$, is a convex set.

Assume (u_1, t_1) and (u_2, t_2) are two points in $\text{epi}(f)$. We want to show that for any λ between 0 and 1, the point $(\lambda u_1 + (1-\lambda)u_2, \lambda t_1 + (1-\lambda)t_2)$ is also in $\text{epi}(f)$. Since (u_1, t_1) is in $\text{epi}(f)$, we have $f(u_1) \leq t_1$. Similarly, (u_2, t_2) being in $\text{epi}(f)$ implies $f(u_2) \leq t_2$.

Using the convexity of f , we have: $f(\lambda u_1 + (1-\lambda)u_2) \leq \lambda f(u_1) + (1-\lambda)f(u_2)$ (convexity property) $\leq \lambda t_1 + (1-\lambda)t_2$ (since $f(u_1) \leq t_1$ and $f(u_2) \leq t_2$)

Therefore, $(\lambda u_1 + (1-\lambda)u_2, \lambda t_1 + (1-\lambda)t_2)$ is in $\text{epi}(f)$, and this shows that $\text{epi}(f)$ is a convex set.

Now, let's prove the other direction. Assume $\text{epi}(f)$ is a convex set, and we want to show that f is a convex function.

Consider any two points u_1 and u_2 in C , and let λ be a scalar between 0 and 1. We need to show that $f(\lambda u_1 + (1-\lambda)u_2) \leq \lambda f(u_1) + (1-\lambda)f(u_2)$.

To do this, we consider the points $(u_1, f(u_1))$ and $(u_2, f(u_2))$ in the epigraph of f . Since $\text{epi}(f)$ is convex, the point $(\lambda u_1 + (1-\lambda)u_2, \lambda f(u_1) + (1-\lambda)f(u_2))$ must also be in $\text{epi}(f)$.

This implies that $f(\lambda u_1 + (1-\lambda)u_2) \leq \lambda f(u_1) + (1-\lambda)f(u_2)$, as desired.

Therefore, we have shown that f is convex if and only if its epigraph, $\text{epi}(f)$, is a convex set.

1.3.4

We will claim that the set F_i (when I belongs to I) is blocked above from the previous Q we know that function is convex if and only if epigraph is a convex set and also $\text{epi}(F) = \{(u, t) | F(u) \leq t\}$ therefore to every I belongs to I : $F_i(u) \leq t$ therefore we get:

$\text{epi}(F) = \{(w, b) \mid \text{Insertion}(F_i) \text{ and because the insertion of convex sets define convex set we conclude that } F \text{ is a convex function.}\}$

1.4 Sub-gradients for Soft-SVM Objective

Based on Lecture 11 and Recitations 2,12

The Soft-SVM objective, though convex, is not differentiable in all of its domain due to the use of the hinge-loss. Therefore, to implement a sub-gradient descent solver for this problem we must first describe sub-gradients of the objective.

5. Given $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{\pm 1\}$. Show that the hinge loss is convex in \mathbf{w}, b . That is, define

$$f(\mathbf{w}, b) := \ell_{\mathbf{x}, y}^{\text{hinge}}(\mathbf{w}, b) = \max(0, 1 - y(\mathbf{x}^\top \mathbf{w} + b))$$

and show that f is convex in \mathbf{w}, b .

6. Deduce some sub-gradient of the hinge loss function $g \in \partial \ell_{\mathbf{x}, y}^{\text{hinge}}(\mathbf{w}, b)$.
 7. Let $f_1, \dots, f_m : \mathbb{R}^d \rightarrow \mathbb{R}$ be a set of convex functions and $\mathbf{g}_k \in \partial f_k(\mathbf{x})$ for all $k \in [m]$ be sub-gradients of these functions. Define $f : \mathbb{R}^d \rightarrow \mathbb{R}$ by $f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x})$. Show that $\sum_k \mathbf{g}_k \in \partial \sum_k f_k(\mathbf{x})$.

8. Let $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \subseteq \mathbb{R}^d \times \{\pm 1\}$ be a sample and define $f : \mathbb{R}^d \rightarrow \mathbb{R}$ by:

$$f(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \ell_{\mathbf{x}_i, y_i}^{\text{hinge}}(\mathbf{w}, b) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Find a sub-gradient of f for any \mathbf{w} .

1.4.5

To prove that a function $f: C \rightarrow \mathbb{R}$ is convex if and only if its epigraph, $\text{epi}(f)$, is a convex set, we can observe that f is a maximum function between two linear functions in w . We need to establish two proofs:

- A linear function is convex.
- The maximum of two convex functions is convex.

Two things we already showed

By demonstrating these two properties, we can conclude that if f is convex, then $\text{epi}(f)$ is convex, and vice versa.

1.4.6

	{	hinge
	{	$(-y\mathbf{x}, -y)$ $(w, b) \neq 0$
	{	\mathbf{x}, y
	{	
	G = {	hinge
	{	0 $(w, b) = 0$
	{	\mathbf{x}, y

[Type here]

1.4.7

By definition of sub-gradient we get that for all i belong to $[m]$, y belong to \mathbb{R}^d :

$$F_i(y) \geq F_i(x) + \langle G_i, y-x \rangle$$

Therefore if we sum m inequality of those we get (i is the index which run on every sigma at this line from 1 to m) $\sum F_i(y) \geq \sum F_i(x) + \langle G, y-x \rangle$

And from the definition of sub-gradient we get that (k is the index which run on each sigma this line) $\sum G_k$ belong to $\partial \sum F_k(x)$ as was needed

1.4.8

We define:

$$G = \begin{cases} \{ (-y_i, -y_i) & \text{hinge} \\ \{ & \text{if } (w, b) \neq 0 \\ \{ & x_i, y_i \\ \{ & \text{hinge} \\ \{ 0 & \text{if } (w, b) = 0 \\ \{ & x_i, y_i \end{cases}$$

(the index at each \sum gonna be i which run from 1 to m)

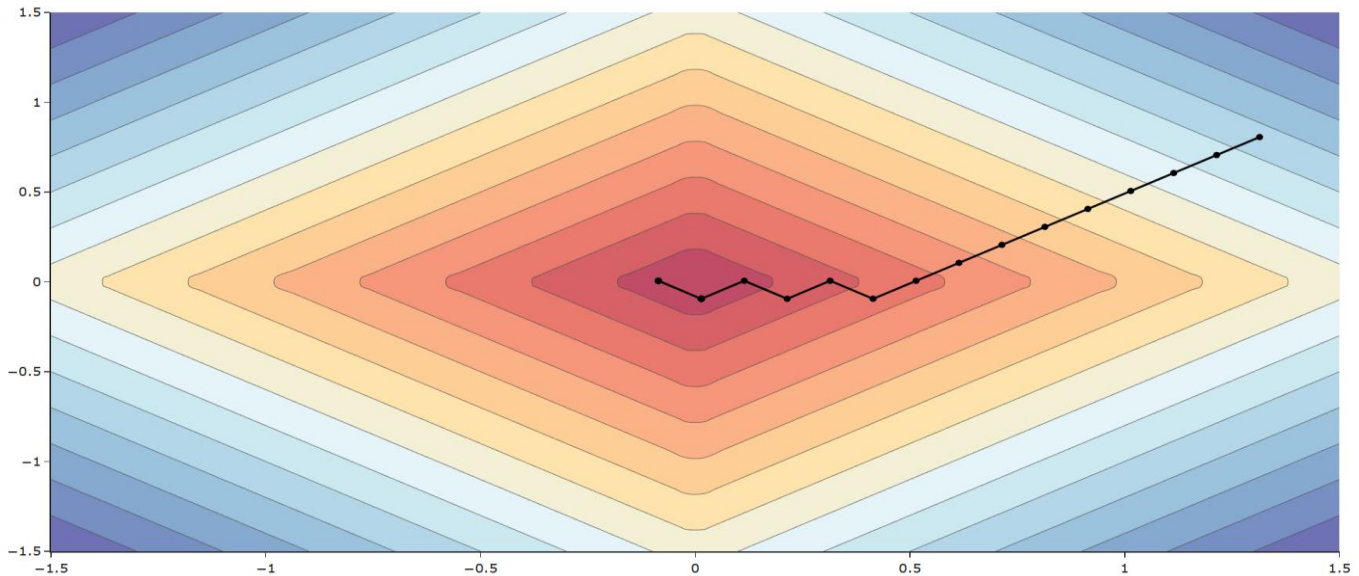
and we get that: $(1/m) \sum G_i + \lambda(w, 0)$ belong to $\partial (1/m) \sum \text{hinge}(w, b) + (\lambda/2) * ||w||^2$

x_i, y_i

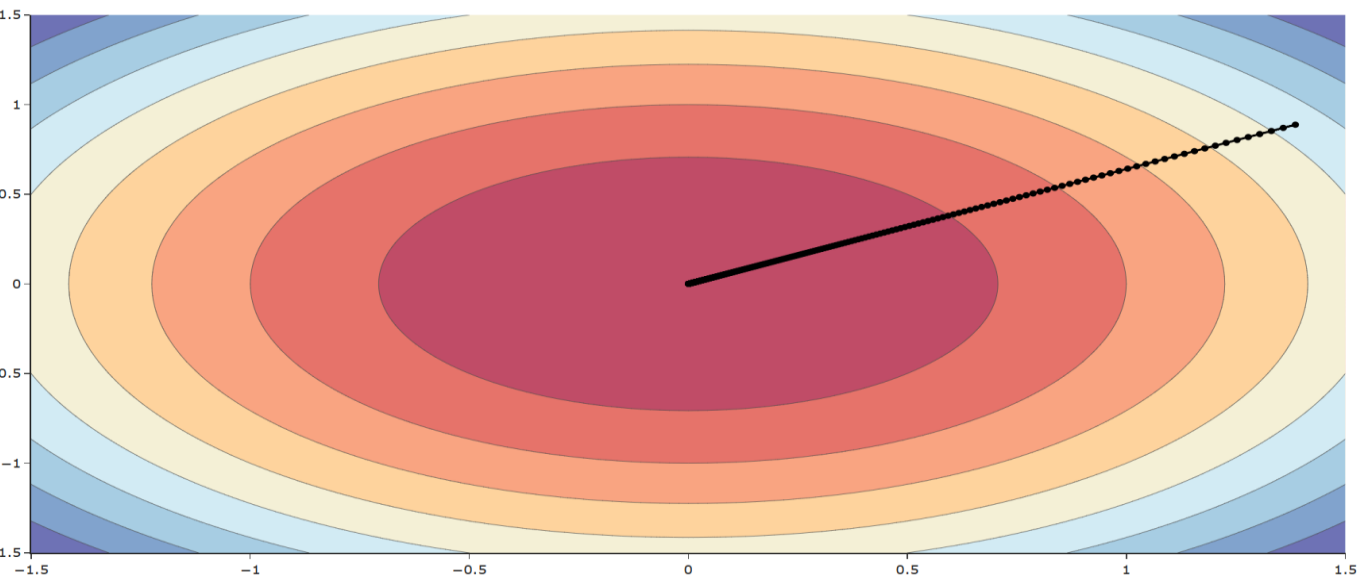
PARTICAL PART

1. Plot the descent path for each of the settings described above (you can use the `plot_descent_path`). Add below the plots for $\eta = 0.01$ and explain the differences seen between the L1 and L2 modules.

GD Descent Path L1 module has eta 0.1



GD Descent Path L2 module has eta 0.01



כפי שלמדנו נורמה 1 יותר מאפסת פיצרים לעומת נורמה 2 שמאפסת פיצרים. ניתן לראות שבגרף של אל 1 הפונקציה מתאפסת כלומר מגיע לאחד הצירים כי חיפשנו שפיץ שיאפס הרבה דבלינו לעומת אל 2 שבה נחפש את נקודת המפגש עבור לוס מינימלי לכן אל 2 מקווצת.

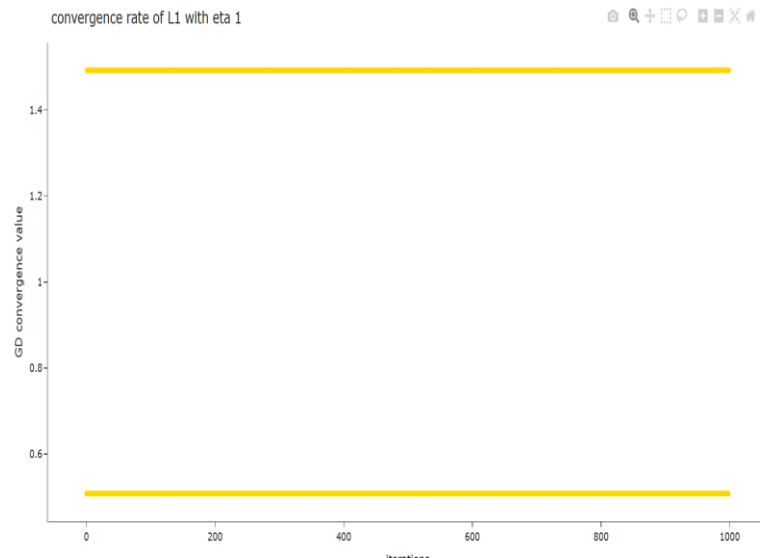
[Type here]

- modules.
- Describe two phenomena that can be seen in the descent path of the ℓ_1 objective when using GD and a fixed learning rate.

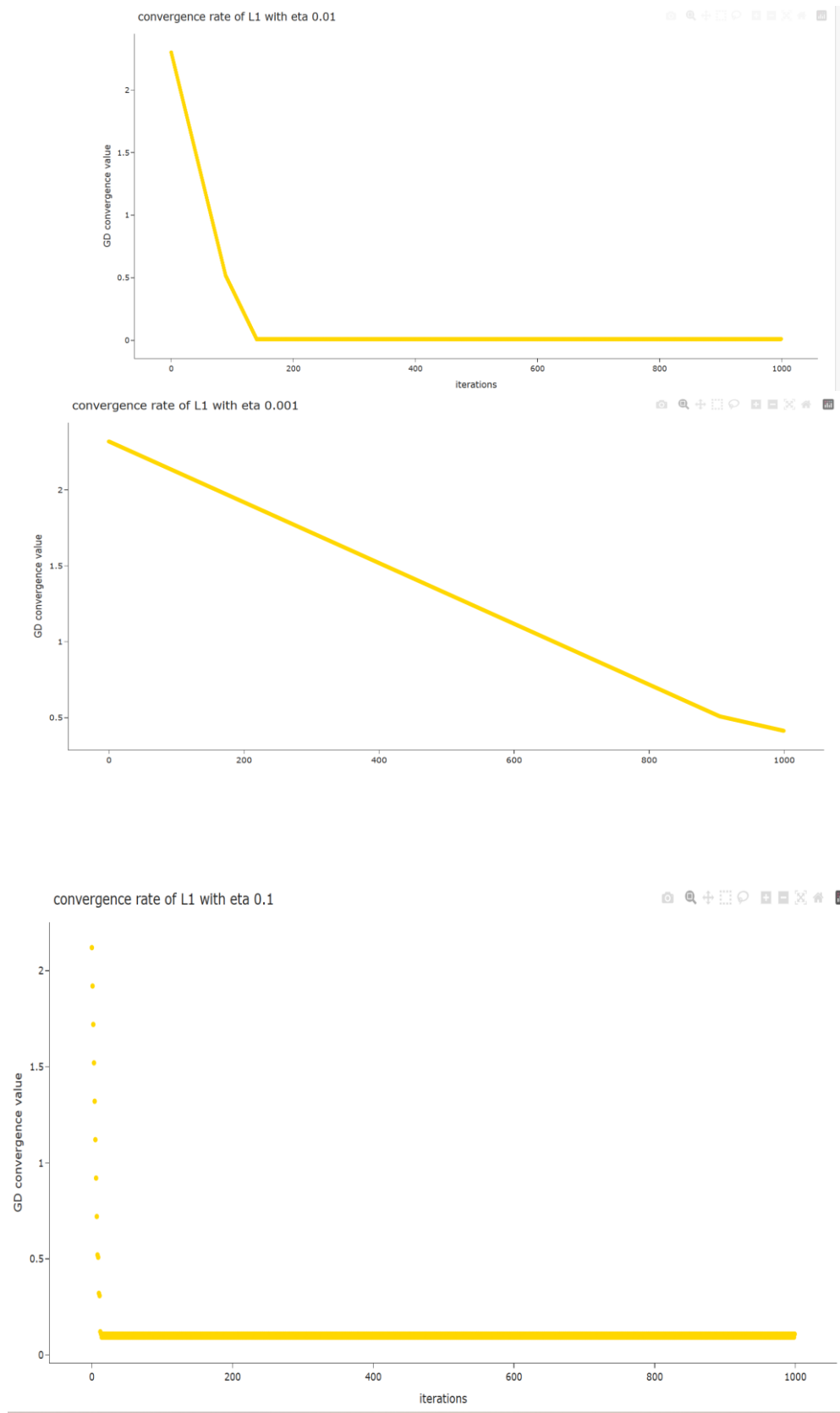
ראשית ניתן להבחין כי אל 1 מנסה לאפס, כלומר להגיע לציר האיקס כמה שיותר מהר
בנוסף לאחר ההגעה לציר האיקס נבחין בקפיצת של זיגזגים סביבו מעל ומתחת

- For each of the modules, plot the convergence rate (i.e. the norm as a function of the GD iteration) for all specified learning rates. Explain your results

L1



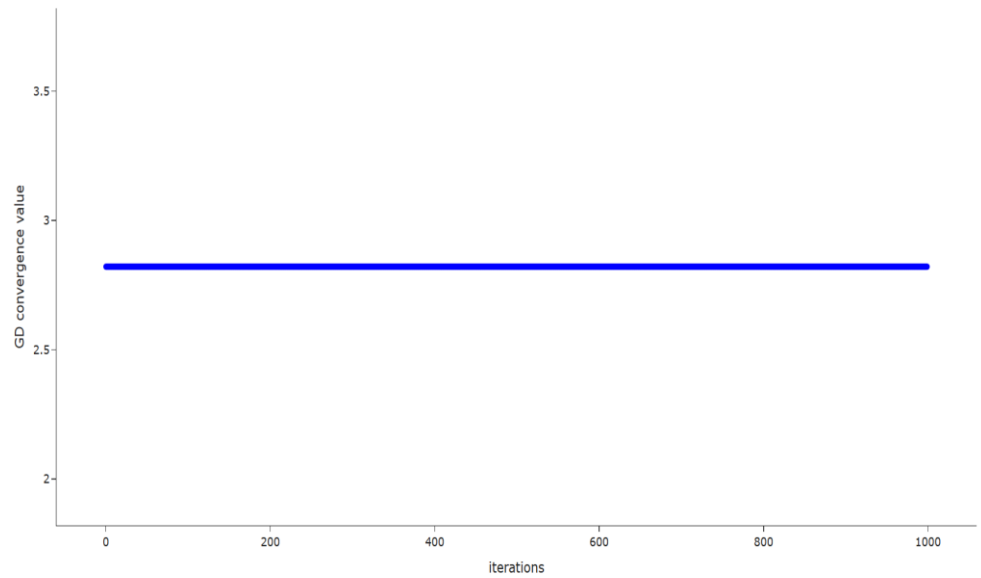
[Type here]



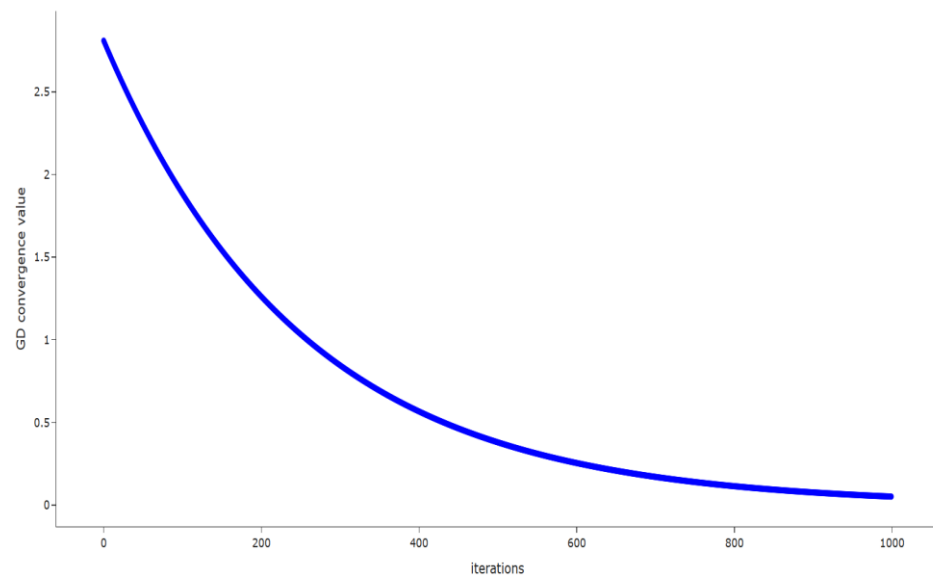
[Type here]

:L2

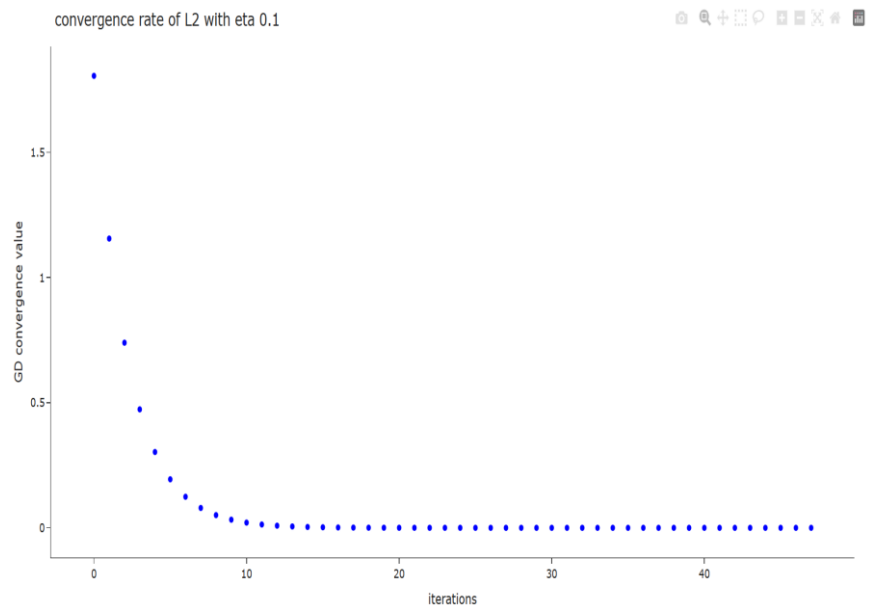
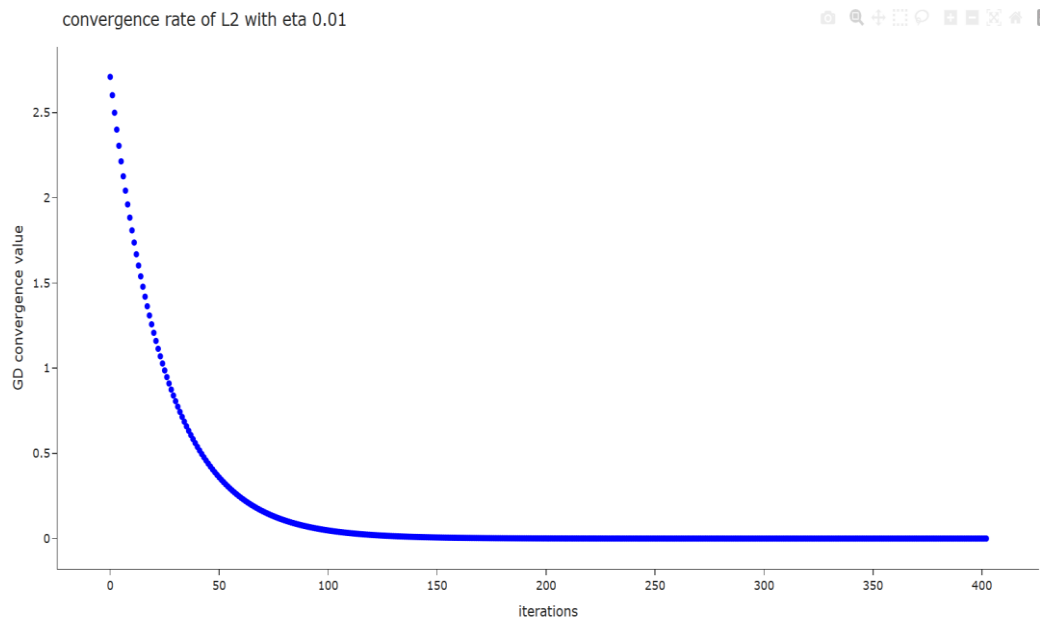
convergence rate of L2 with eta 1



convergence rate of L2 with eta 0.001



[Type here]



נסתכל על הגרפים, קל לראות כי מהגדרת אל 1 ואל 2 עבור אל 1 נקבל פונקציה עם הרבה נקודות אי רציפות שמתאפס מהר לעומת אל 2 שיותר רציפה ומתאפסת לאט יותר כלומר מתכווצת יחד לאט לאט

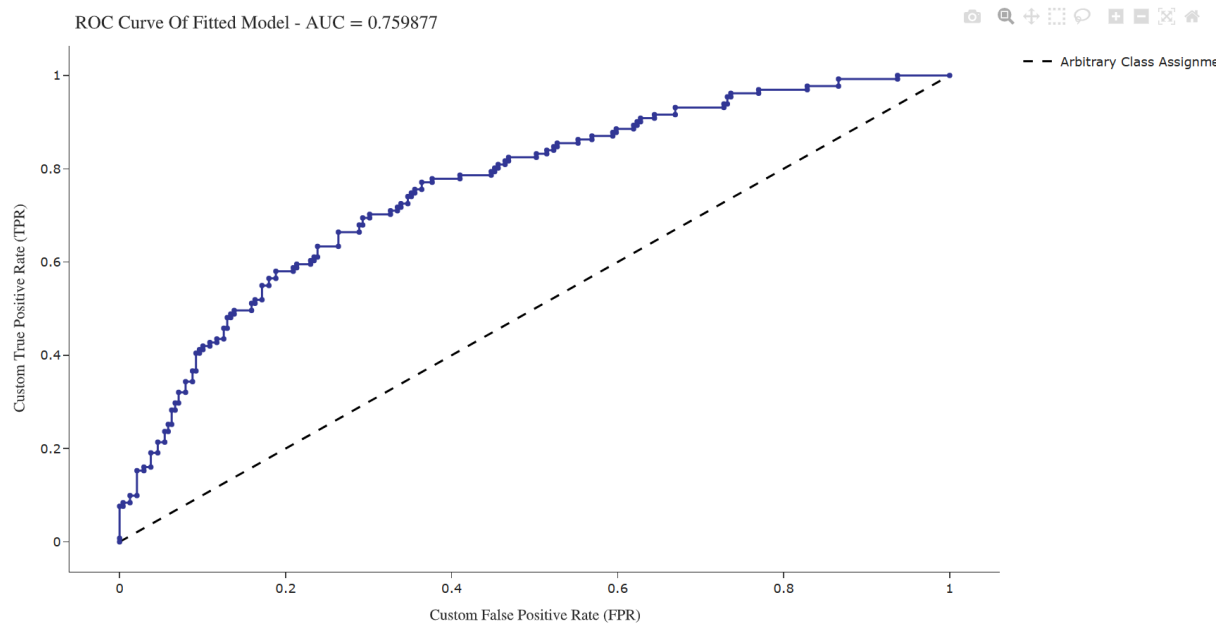
4. What is the lowest loss achieved when minimizing each of the modules? Explain the differences

```
best loss L1 module is: 0.011880380446586989 with eta 0.01
best loss L2 module is: 1.4029519498344255e-09 with eta 0.1
```

מכך שאלו קופץ בין הערכים אנו מעט רחוקים יותר מ0 לעומת אל 2 שבכל צעד מקוץ את הפונקציה ויורד מונוטונית ל0 נקבל ערך ממש קטן

Then, load the South Africa Heart Disease dataset (SAheart.data), split it to train- and test sets (80% train) and answer the following questions:

8. Using your implementation, fit a logistic regression model over the data. Use the `predict_proba` to plot an ROC curve. You can use sklearn's `metrics.roc_curve` function and the code provided in Lab 04.



9. Which value of α achieves the optimal ROC value according to the criterion below. Using this value of α^* what is the model's test error?

$$\alpha^* = \operatorname{argmax}_{\alpha} \{ \operatorname{TPR}_{\alpha} - \operatorname{FPR}_{\alpha} \}$$

```
The best alpha is: 0.32
Model's test error: 0.33695652173913043
```

[Type here]

10. Fit an ℓ_1 -regularized logistic regression by passing `penalty="l1"` when instantiating a logistic regression estimator
- Set $\alpha = 0.5$
 - Use your previously implemented cross-validation procedure to choose λ
 - After selecting λ repeat fitting with the chosen λ and $\alpha = 0.5$ over the entire train portion.

For values of What value of λ was selected and what is the model's test error?

11. Repeat question 10 for ℓ_2 regularized logistic regression. What value of λ was selected and what is the model's test error?

Best lamda L1 is 0.02 and its loss is 0.27

Best lamda L2 is 0.01 and its loss is 0.29