

Principles of Machine Learning Exercises

Victor Verreet victor.verreet@cs.kuleuven.be
Laurens Devos laurens.devos@cs.kuleuven.be

Fall, 2021

Exercise Session 3: Optimization and ILP

3.1 Huber Loss Gradient

Huber loss is defined as

$$\mathcal{L}_{H,\delta}(y, \hat{y}) = \begin{cases} \frac{1}{2}(y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq \delta, \\ \delta|y - \hat{y}| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

- Draw this curve.
- Huber is made up of two other loss functions we've seen in the course. Why would you use Huber instead of the other losses?
- Compute the gradient of the Huber loss. Is this function C^1 continuous? Is it C^2 continuous? Do the answers to the previous two questions limit the possible optimization methods?
- Huber loss has a parameter δ . Which value would you choose for this parameter?

3.2 Squared Loss with Squared Regularization

- Write down the loss function $\mathcal{L}(y, Xw)$ for a linear regression problem with L2 regularization with input matrix X containing the instances in its rows, a column vector with weights w , and a column vector y with target values.
- Compute the derivative of this loss function.¹

3.3 Non-Differentiable Functions: Subgradients

The following definition for convexity was given in the lecture: for any $x, y \in \text{dom}(f)$ and $\alpha \in [0, 1]$:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y).$$

If a function f is differentiable, f is convex if:²

Fix x

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) \quad \text{for all } x, y \in \text{dom}(f).$$

- Visually show what these definitions mean for $f(x) = x^2$. Explain how what you have drawn relates to the definition.

¹The *Dive Into Deep Learning* book has a great chapter on multivariate calculus.

²A great resource for convex optimization is: <http://www.stat.cmu.edu/~ryantibs/convexopt>.

- b. A **subgradient** of a convex function (also non-differentiable) at a value x is any vector $g \in \mathbb{R}^n$ such that:

$$f(y) \geq f(x) + g^T(y - x) \quad \text{for all } y \in \text{dom}(f).$$

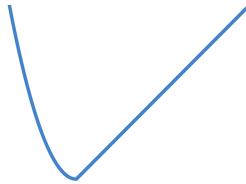
The **subgradient set** or **subdifferential** is defined as all vectors g for which the above holds:

$$\partial f(x) = \{g \mid g \text{ is a subgradient of } f \text{ at } x\}.$$

Find the subgradient set $\partial f(x)$ for all $x \in \mathbb{R}$ of the function:

$$f(x) = \begin{cases} x^2, & x \leq 0, \\ x, & x > 0. \end{cases}$$

$$f(y) \geq f(x) + g^T(y - x)$$



$$\begin{aligned} g &= 0 \\ \Rightarrow f(y) &\geq f(x) \end{aligned}$$

A convex function is minimized at a value x^* if $0 \in \partial f(x^*)$. Can you verify this in the above example?

3.4 First-Order Methods Using Parabolas

We've seen that a second-order approximation of a differentiable function can be used to optimize a function (e.g. Newton-Raphson). The general format of such second-order approximation in a is the following:

$$f(x) \approx f(a) + \nabla f(a)^T(x - a) + \frac{1}{2}(x - a)^T \mathbf{H}(a)(x - a),$$

where $\mathbf{H}(a)$ is the Hessian of f evaluated in a . We can turn this into a first-order approximation by replacing the Hessian by the scaled identity matrix $\frac{1}{\rho}\mathbf{I}$, i.e., we still use a parabolic approximation, but we limit the family of possible approximating functions to the parabolas for which the coefficients for the second-order terms are fixed. You can view ρ as a scaling factor of the parabola. We get:

$$\begin{aligned} f(x) &\approx f(a) + \nabla f(a)^T(x - a) + \frac{1}{2\rho}(x - a)^T \mathbf{I}(x - a) \\ &= f(a) + \nabla f(a)^T(x - a) + \frac{1}{2\rho} \|x - a\|_2^2. \end{aligned}$$

Show that:

$$\arg \min_x f(a) + \nabla f(a)^T(x - a) + \frac{1}{2\rho} \|x - a\|_2^2$$

is equivalent to

$$\arg \min_x \frac{1}{2\rho} \|(x - (a - \rho \nabla f(a)))\|_2^2.$$

Note that the above formula can be solved exactly: simply take x equal to $a - \rho \nabla f(a)$, which corresponds to the **update step of gradient descent**, with ρ the step size or learning rate.

3.5 Proximal Methods: Finding Weights for LASSO

Proximal gradient descent is used to optimize convex functions of the form $f(x) = g(x) + h(x)$, where f is non-differentiable, but it can be decomposed into a convex differentiable g and a convex non-differentiable h . We cannot apply the gradient update step from the previous exercise to $f(x)$ in its

entirely because it is non-differentiable. However, we can use it for $g(x)$. Consider the following update step, where $g(x)$ is replaced by its quadratic approximation $\tilde{g}(x)$ and $h(x)$ is left unchanged:

$$\begin{aligned} x^{(k+1)} &= \arg \min_z \tilde{g}(z) + h(z) \\ &= \arg \min_z g(x^{(k)}) + \nabla g(x^{(k)})^T (z - x^{(k)}) + \frac{1}{2\rho} \|z - x^{(k)}\|_2^2 + h(z) \\ &= \arg \min_z \underbrace{\frac{1}{2\rho} \|z - (x^{(k)} - \rho \nabla g(x^{(k)}))\|_2^2}_{(1)} + \underbrace{h(z)}_{(2)} \end{aligned} \quad (1)$$

We can interpret (1) as: choose an update z as close to the gradient update $x - \rho \nabla g(x)$, and (2) as: but also make $h(z)$ small. Proximal gradient descent uses a proximal mapping to optimize functions of this form. We won't be looking at the general formulation, rather, we will look at an example.

Proximal gradient descent can be used to optimize the LASSO criterion for linear regression. Linear regression with L1 regularization (LASSO) is defined as:

$$\begin{aligned} &\frac{1}{2} \sum_i (y_i - \mathbf{X}_{i*} w)^2 + \lambda \|w\|_1 \\ &= \frac{1}{2} \|y - \mathbf{X} w\|_2^2 + \lambda \|w\|_1. \end{aligned} \quad (2)$$

The regularization term is not differentiable at zero, so standard gradient descent models aren't applicable. Luckily, the problem can be decomposed into $f = g + h$, where g is the convex differentiable part $\sum_i (y_i - \mathbf{X}_{i*} w)^2$, and h is the convex, non-differentiable part $\lambda \|w\|_1$.

We will use the proximal mapping of the L1 norm. But what exactly is a proximal mapping? The general definition of a *proximal mapping* or *proximal operator* of a convex function ϕ is:

$$\text{prox}_\phi(v) = \arg \min_z \left[\phi(z) + \frac{1}{2} \|z - v\|_2^2 \right]$$

Complete books have been written on the subject of proximal methods, so the full details are beyond the scope of this course. What you need to know is that $\text{prox}_\phi(v)$ is well-defined: it has a unique solution because $\|\cdot\|_2$ is strictly convex and we required ϕ to be convex. The sum is therefore also strictly convex, so a unique solution exists. Intuitively, you can view the proximal mapping as a projection to a point in the neighborhood of v that is (1) close to v and (2) minimizes ϕ .

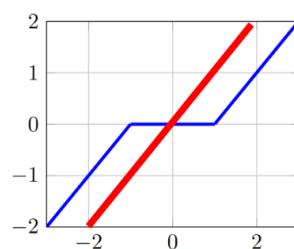
Replacing ϕ in the general definition of a proximal mapping with the L1 norm gives (with weight $\lambda > 0$):

$$\text{prox}_{\lambda \|\cdot\|_1}(v) = \arg \min_z \left[\lambda \|z\|_1 + \frac{1}{2} \|z - v\|_2^2 \right] \quad (3)$$

a. Show that the above is minimized for:

$$\begin{aligned} z^* &= S_\lambda(v) \\ [S_\lambda(v)]_i &= \begin{cases} v_i - \lambda & \text{if } v_i > \lambda \\ 0 & \text{if } -\lambda \leq v_i \leq \lambda \\ v_i + \lambda & \text{if } v_i < -\lambda \end{cases} \end{aligned} \quad (4)$$

Remember that a convex function ρ is minimized for an input z^* when $0 \in \partial\rho(z^*)$. The S_λ function is called the **soft-thresholding operator**:



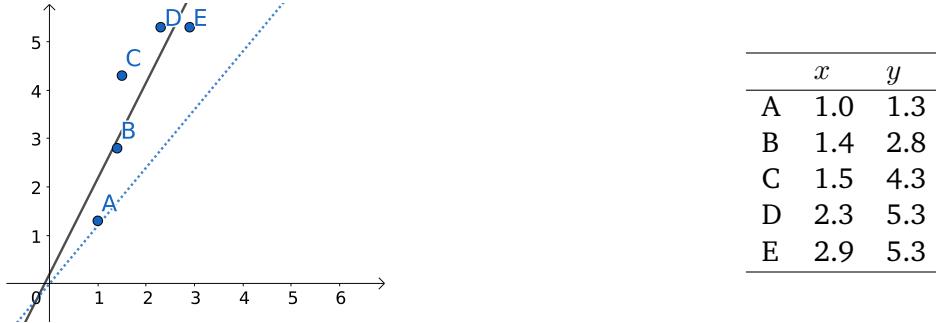
- b. Show that, when \mathbf{X} is taken to be the identity matrix \mathbf{I} , LASSO in Equation 2 is equivalent to $\text{prox}_{\lambda \|\cdot\|_1}$ from Equation 3.
- c. Rewrite the update step in Equation 1 for LASSO (Equation 2). Show that this is equivalent to:

$$\begin{aligned} w^{(k+1)} &= S_{\lambda\rho} \left(w^{(k)} - \rho \nabla g(w^{(k)}) \right) \\ &= S_{\lambda\rho} \left(w^{(k)} + \rho \mathbf{X}^T (y - \mathbf{X}w^{(k)}) \right), \end{aligned}$$

where $\nabla g(w) = -\mathbf{X}^T(y - \mathbf{X}w)$ (see previous exercise). This is the updating step of the **iterative soft-thresholding algorithm (ISTA)**.

- d. Given the dataset below, $\lambda = 0.2$, learning rate $\rho = 0.01$, and $w^{(0)} = [0.0; 0.5]$, compute an update step using iterative soft thresholding. Compute the objective function before and after applying the soft-thresholding operator. Is the objective function minimized?

Feel free to use Python, Matlab, Julia, etc. for this exercise. This also allows you to execute multiple iterations.



3.6 Theta Subsumption

Order these clauses according to theta subsumption. Which clauses are reduced? What are the reductions of the non-reduced clauses?

1. $f(a, b) \leftarrow m(a), m(b), p(a, b).$
2. $f(X, Y) \leftarrow m(X), p(X, Y).$
3. $f(X, X) \leftarrow m(X), p(X, Y), p(X, Z).$
4. $f(X, X) \leftarrow m(X), p(X, X), m(X).$
5. $f(X, X) \leftarrow m(X), p(X, X).$
6. $f(X, X) \leftarrow m(X), p(X, Y).$

3.7 Least General Generalization

Compute the LGG of the following literals and clauses.

- a. $m(a, c(a, \text{nil}))$ $m(x, c(y, \text{nil}))$
 $m(b, c(c, \text{nil}))$
- b. $\text{add}(s(s(0)), s(s(0)), s(s(s(s(0))))))$ $\text{add}(s(x), n(x), s(z))$
 $\text{add}(s(0), s(0), s(s(0)))$
- c. $m(X, c(X, Y)) \leftarrow c(X), \text{list}(Y).$
 $m(X, c(Y, Z)) \leftarrow c(X), c(Y), \text{list}(Z), m(X, Z).$

3.8 Subsumption and Implication

Show that if $A \leq_\theta B$ then $A \implies B$ holds for any clauses A and B .

3.9 Cyclic Clauses

Define for every $n > 0$ a clause γ_n as

$$a \iff \bigwedge_{i=0}^{n-1} c(X_i, X_{i+1})$$

where we define $X_n = X_0$. Give the LGG of two clauses γ_n and γ_m for general n and m and prove this result.

3.10 Inverse Resolution

Find an inverse resolution path that derives the clause

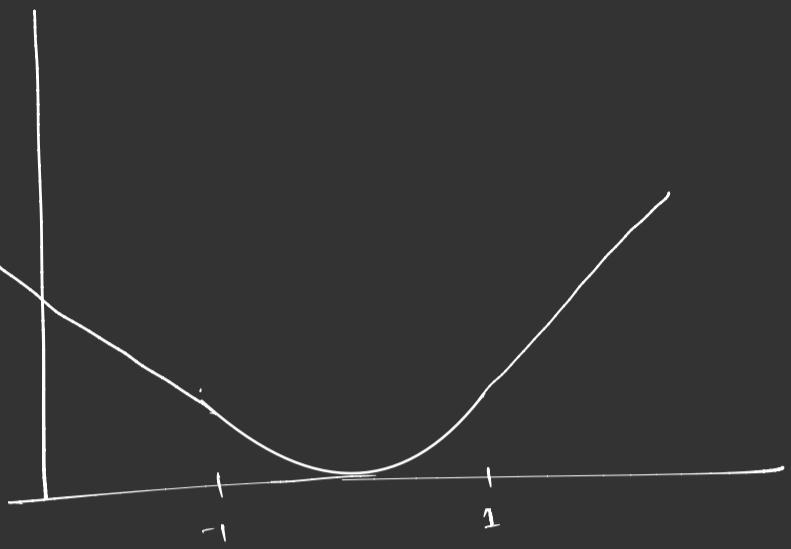
$$\text{daughter}(X, Y) \iff \text{parent}(Y, X), \text{female}(X)$$

from the observed facts

$$\begin{aligned} &\text{daughter(lisa, marge)} \\ &\text{parent(marge, lisa)} \\ &\text{female(lisa)} \end{aligned}$$

$$\textcircled{1} \quad \ell_{H,\delta}(y, \hat{y}) = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{for } |y - \hat{y}| \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

Let $\delta = 1$

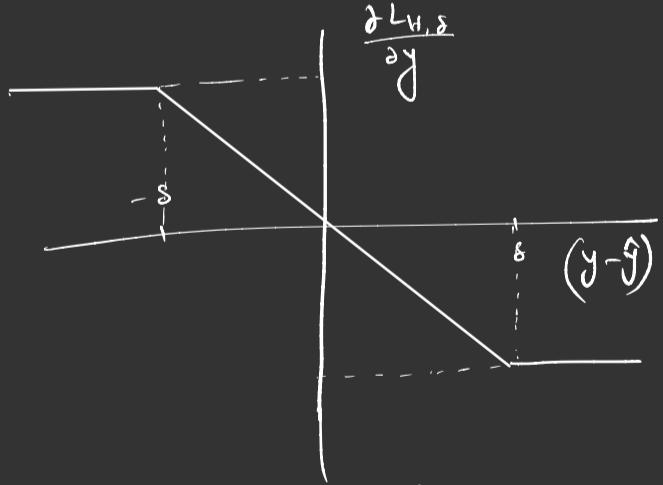


Huber Loss has robustness like L1 and also differentiable at 0.

$$\frac{\partial \ell}{\partial \hat{y}} = -(y - \hat{y}) \quad \text{if } |y - \hat{y}| \leq \delta$$

$$\frac{\partial \ell}{\partial \hat{y}} = \begin{cases} -\delta & \text{if } (y - \hat{y}) > \delta \\ \delta & \text{if } (y - \hat{y}) < -\delta \end{cases}$$

$$\frac{\partial \ell}{\partial \hat{y}} = \begin{cases} \delta & \text{if } (y - \hat{y}) < -\delta \\ -(y - \hat{y}) & \text{if } |y - \hat{y}| \leq \delta \\ -\delta & \text{if } (y - \hat{y}) > \delta \end{cases}$$



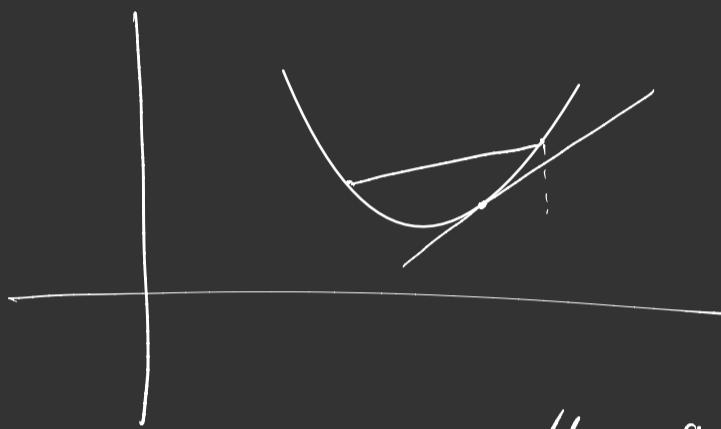
$$\frac{\partial^2 \ell}{\partial \hat{y}^2} = \begin{cases} 0 & \text{if } |y - \hat{y}| > \delta \\ 1 & \text{if } |y - \hat{y}| \leq \delta \\ 0 & \text{otherwise} \end{cases}$$

$$\underline{3.2} \quad \mathcal{L}(y, xw) = \frac{1}{2} \|y - xw\|^2 + \lambda \|w\|^2$$

$$\frac{\partial \mathcal{L}}{\partial w} = \begin{bmatrix} -2 \cdot \frac{1}{2} (y_i - x_i w_i) x_{i1} \\ \vdots \\ -2 \cdot \frac{1}{2} (y_n - x_n w_n) x_{n1} \end{bmatrix} + 2\lambda w$$

$$= x^T (xw - y) + 2\lambda w$$

$$\frac{f(y) - f(x)}{y - x}$$



Eq^n of a Tangent at pt. $x=a$

$$T(x) = f(a) + \nabla f(a)(x-a)$$

$$\forall x \in X$$

$$T(x) \leq f(x)$$

Tangent will lie under the curve.

$$\begin{aligned} f(y) - 2xy + x^2 &\geq 0 \\ \Rightarrow (x-y)^2 &\geq 0 \\ y^2 &\geq x^2 \end{aligned}$$

$$\text{b) } \partial f(x) = \begin{cases} 2x & x < 0 \\ [0, 1] & x=0 \\ 1 & x > 0 \end{cases}$$

3.4

$$g(x) = f(a) + \nabla f(a)^T(x-a) + \frac{1}{2} \rho \|x-a\|_2^2$$

$$\frac{\partial g}{\partial x} = \nabla f(a)^T + \frac{1}{\rho}(x-a) = 0$$

$$\Rightarrow x = a - \rho \nabla f(a)^T$$

$$x = a - \rho \nabla f(a)^T$$

$$\text{Prox}_\phi(v) = \arg_{z \in \mathbb{R}} \left(\phi(z) + \frac{1}{2} \|z-v\|^2 \right)$$

$$r(v) = \text{Prox}_{\lambda \|\cdot\|_1}(v) = \arg_{z \in \mathbb{R}} \left(\lambda \|z\|_1 + \frac{1}{2} \|z-v\|^2 \right)$$

$$\text{If } z > 0, z = v - \lambda, v > \lambda$$

$$z=0 \quad -\lambda \leq v \leq \lambda$$

$$z = v + \lambda \quad v < -\lambda$$

$$\arg_{x_i} \min f(a) + \nabla f(a)^T(x_i - a_i) + \frac{1}{2\rho} \|x_i - a_i\|^2$$

$$\arg_{x_i} \min \sum_i \partial_i f(a) (x_i - a_i) + \frac{1}{2\rho} \sum_i (x_i - a_i)^2$$

$$\arg_{x_i} \min \frac{1}{2\rho} \sum_i \left(2\rho \partial_i f(a) - 2\rho a_i \partial_i f(a) + x_i^2 - 2x_i a_i + a_i^2 \right)$$

$$= \arg_{x_i} \min \frac{1}{2\rho} \sum_i x_i^2 + 2\rho x_i \partial_i f(a) - 2\rho a_i \partial_i f(a) + a_i^2 - 2x_i a_i + \left(\rho \partial_i f(a) \right)^2 - \left(\rho \partial_i f(a) \right)^2$$

$$= \arg_{x_i} \min \frac{1}{2\rho} \sum_i \left\{ x_i^2 - 2x_i (a_i - \rho \partial_i f(a)) + (a_i - \rho \partial_i f(a))^2 \right\}$$

$$= \arg_{x_i} \min \frac{1}{2\rho} \sum_i \left[x_i - (a_i - \rho \partial_i f(a))^2 \right]^2$$

$$= \arg_{x_i} \min \frac{1}{2\rho} \|x_i - (a_i - \rho \nabla f(a))^2\|_2^2$$

$$x^{k+1} = \arg_{z \in \mathbb{R}} \tilde{g}(z) + h(z)$$

$$= \arg_{z \in \mathbb{R}} \left(g(x^k) + \nabla g(x^k)^T(z - x^k) + \frac{1}{2\rho} \|z - x^k\|_2^2 \right) + h(z)$$

$$= \arg_{z \in \mathbb{R}} \left\| z - (x^k - \rho \nabla g(x^k))^2 \right\|_2^2 + h(z)$$

$$f = g + h$$

$$p(v) = \arg_{z \in \mathbb{R}} \left(\lambda \|z\|_1 + \frac{1}{2} \|z - v\|_2^2 \right)$$

$$\partial \phi(z) = \begin{cases} \lambda s + (z - v) & |s| \in \partial(\|z\|_1) \end{cases}$$

$$s = \begin{cases} -1 & \text{if } z < 0 \\ [-1, 1] & \text{if } z = 0 \\ 1 & \text{if } z > 0 \end{cases}$$

$$\text{If } v > \lambda > 0, v - \lambda s > 0$$

$$\lambda s + z - v = 0$$

$$\Rightarrow z = v - \lambda s$$

-3 < -2

-3 + 2

$$s = 1, z = v - \lambda$$

$$\text{If } v < -\lambda, z = v - \lambda s < -\lambda - \lambda s < 0$$

$$s = -1$$

$$z = v + \lambda$$

$$\text{If } -\lambda \leq v \leq \lambda$$

$$v - \lambda s = 0 \text{ for } s = \frac{v}{\lambda} \in [-1, 1]$$

$$\therefore z = 0$$

$$\frac{1}{2} \|\mathbf{y} - \mathbf{x}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

$$= \frac{1}{2} \|\mathbf{y} - \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

$$\mathbf{w}^{k+1} = \arg\min \left(\lambda \|\mathbf{w}\|_1 + \frac{1}{2} \|\mathbf{y} - \mathbf{w}\|_2^2 \right)$$

$$\mathbf{w}^{k+1} = \underset{\lambda, p}{\operatorname{Solve}} \left(\mathbf{w}^k - p \nabla g(\mathbf{w}^k) \right)$$

$$= \arg\min_{\mathbf{z}} \left(\lambda \|\mathbf{z}\|_1 + \frac{1}{2} \|\mathbf{y} - (\mathbf{w}^k - p \nabla g(\mathbf{w}^k))\|_2^2 \right)$$

$$\nabla g(\mathbf{w}) = -\mathbf{x}^\top (\mathbf{y} - \mathbf{x}\mathbf{w})$$

$$\mathbf{w}_{0,0} = \begin{pmatrix} 0 & 0 \\ 0 & 5 \end{pmatrix}, \quad \lambda = 0.2, \quad p = 0.01$$

$$\textcircled{1} \quad f(a, b) \leftarrow m(a), m(b), p(a, b)$$

$$2 \leq 1$$

