

Section 4: The LASSO and gradient descent

How to use a gradient descent to solve the LASSO?

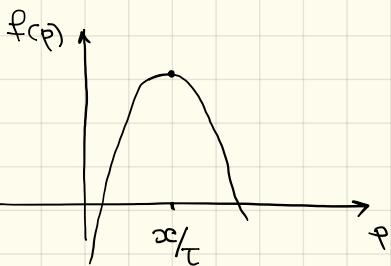
The main difficulty to overcome is non-differentiability of the absolute value function. The trick is to replace the absolute value with a differentiable function.

$$|x| = \max_{p \in [-1, 1]} xp$$

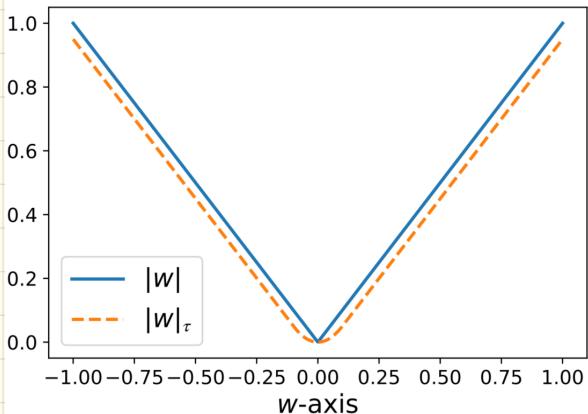
$$|x|_\tau := \max_{p \in [-1, 1]} xp - \frac{\tau}{2} p^2, \quad \tau > 0$$

↑ small perturbation

$$|x|_\tau - ?$$



- $|x| \leq \tau \Rightarrow |x|_\tau = \frac{x^2}{2\tau}$
- $|x| > \tau \Rightarrow |x|_\tau = |x| - \frac{\tau^2}{2}$



$$\tau = \frac{1}{10}$$

Properties of the smoothing:

1. The smooth absolute value is differentiable.

$$\frac{d}{dw} |w|_\tau = \begin{cases} 1, & w > \tau \\ \frac{1}{\tau} w, & |w| \leq \tau \\ -1, & w < -\tau \end{cases}$$

2. The smooth absolute value is close to abs.

$$|w_\tau| \leq |w| \leq |w_\tau| + \frac{\tau}{2}.$$

$$w^* = \arg \min_w \frac{1}{2} \| \Phi(X)w - y \|^2 + \alpha H_\tau(w)$$

$$H_\tau(w) = \sum_{i=0}^d |w_i|_\tau.$$

Can we use normal equation? - No, the system is not linear anymore.

Can we use gradient descent? - Yes, with minor upgrade.

$$E_\tau(w) = \frac{1}{2} \|Xw - y\|^2 + \alpha H_\tau(w)$$

$$w^{k+1} = w^k - \tau [x^\top (Xw^k - y) + \alpha \nabla H_\tau(w^k)]$$

! in practice

forward-splitting

$$w^{k+\frac{1}{2}} = w^k - \frac{\tau}{\alpha} x^\top (Xw^k - y)$$

$$w^k = w^{k+\frac{1}{2}} - \tau \nabla H_\tau(w^{k+\frac{1}{2}})$$

Remark: $w - \tau \nabla |w_\tau| = \begin{cases} w - \tau, & \geq \tau \\ 0, & \\ w + \tau, & w \leq -\tau \end{cases} =: \text{soft}_\tau(w)$
 soft-tresholding

$$w_j^{k+1} = \text{soft}_\tau \left[\left(w^k - \frac{\tau}{2} X^\top (Xw^k - y) \right)_j \right]$$

This is known as ISTA (iterative soft-thresholding algorithm) and is a special case of proximal gradient descent

$$w^{k+1} = (I + \tau \nabla R)^{-1} (w^k - \tau \nabla L(w^k))$$

$\xrightarrow[\text{proximal map}]{} (I + \tau \nabla R)^{-1}(z) := \arg \min_w \left\{ \frac{1}{2} \|w - z\|^2 + \tau R(w) \right\}$

$$L = \frac{1}{2\alpha} \|Xw - y\|^2$$

$$R = \tau \| \cdot \|_1$$

Section 5 Proximal gradient descent

Recap: The LASSO problem consists of

$$W_\alpha = \arg \min_w \left\{ \frac{1}{2} \|Xw - y\|^2 + \alpha H_T(w) \right\}$$

where $H_T(w) = \sum_{j=0}^d |w_j|_T$, and $|x|_T = \begin{cases} |x| - \frac{T}{2}, & |x| > T \\ \frac{1}{2T} |x|^2, & |x| \leq T \end{cases}$

Gradient descent in this case has a form

$$W^{k+1} = W^k - \tau \left[X^\top (Xw^k - y) + \alpha \nabla H_T(w^k) \right]$$

Alternative: proximal gradient descent

Suppose we want to minimise

$$E(w) = L(w) + R(w)$$

where

1. L is convex, continuously differentiable
2. R is proper, convex, lower-semi continuous and has a simple proximal map, i.e.

proximal map

$$\text{PROX}_R(z) := \arg \min_x \left[\frac{1}{2} \|x - z\|^2 + R(x) \right]$$

is easy to compute

↑ distance to z

Def. The proximal gradient descent method is an iterative procedure of the form

$$W^{k+1} = \text{PROX}_R \left(W^{(k)} - \tau \nabla L(W^{(k)}) \right)$$

for the energy $E(w) = L(w) + R(w)$

If additionally to the above R is differentiable

$$\nabla \left(\frac{1}{2} \|x - z\|^2 + \tau R(x) \right) = x - z + \tau \nabla R(x) = 0$$

$$w^{k+1} + \tau \nabla R(w^{k+1}) - (w^k - \tau \nabla L(w^k)) = 0$$

$$w^{k+1} = w^k - \tau (\nabla L(w^k) + \nabla R(w^{k+1}))$$

↑ implicit-explicit gradient method

$$\text{Example: } R(w) = \frac{1}{2} \|w\|^2$$

$$\text{prox}_{\tau R}(z) = \arg \min_x \left\{ \frac{1}{2} \|x - z\|^2 + \underbrace{\frac{\tau}{2} \|x\|^2}_{E(x)} \right\}$$

$$\nabla E(x) = x - z + \tau x = 0$$

$$\hat{x} = \frac{1}{1+\tau} \cdot z$$

$$(1 + \tau \triangleright \frac{1}{2} \| \cdot \|^2)^{-1}(z) = \frac{z}{1+\tau}$$

$$\text{Example: } R(w) = \begin{cases} 0, & w \in C \\ \infty, & w \notin C \end{cases} \quad \text{↑ convex set}$$

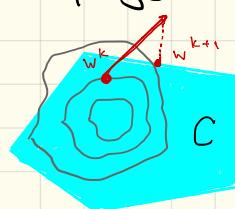
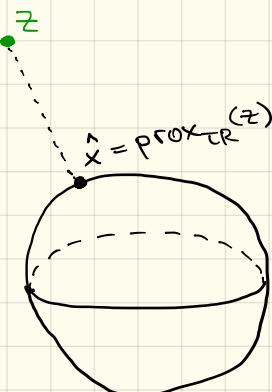
Remark: This corresponds to a minimisation problem with constraint $w \in C$.

$$\text{prox}_{\tau R}(z) = \arg \min_x \left[\frac{1}{2} \|x - z\|^2 + \tau R(x) \right]$$

$$= \arg \min_{x \in C} \frac{1}{2} \|x - z\|^2$$

projection onto convex set

$$w^{k+1} = \text{proj}_C (w^k - \tau \nabla L(w^k))$$



$$\text{Example: } C = \{x \in \mathbb{R} : x \geq 0\}$$

$$\text{proj}_C(z) = \max(0, z) \quad [\text{Exercise}]$$

$$w^{k+1} = \max(0, w^k - \tau \nabla L(w^k))$$

↑ ensures that the iterate stays positive.

What is the motivation of proximal gradient descent?

$$w^{k+1} = \arg \min_w \{L(w) + R(w) + D_J(w, w^k)\}$$

$$\text{where } J = \frac{1}{2\tau} \|w\|^2 - L(w)$$

$$D_J(u, v) = J(u) - J(v) - \langle \nabla J(v), u - v \rangle$$

↑ analogue of $\frac{1}{2} J''(v) \cdot (u - v)^2$ in $d=1$

$$w^{k+1} = \arg \min_w \left\{ L(w) + R(w) + \frac{1}{2\tau} D_{\| \cdot \|_2}(w, w^k) - D_L(w, w^k) \right\}$$

$$= \arg \min_w \left\{ L(w^k) - \langle \nabla L(w^k), w - w^k \rangle + R(w) + \frac{1}{2\tau} \left(\|w\|^2 - \|w^k\|^2 - 2 \langle w^k, w - w^k \rangle \right) \right\}$$

$$= \arg \min_w \left\{ \tau R(w) + \frac{1}{2} \|w - w^k\|^2 - \tau \langle \nabla L(w^k), w - w^k \rangle \right\}$$

$$= \arg \min_w \left\{ \frac{1}{2} \|w - (w^k - \tau \nabla L(w^k))\|^2 + \tau R(w) \right\}$$

How do we use proximal gradient descent to solve the LASSO problem?

$$\hat{w}_2 = \arg \min_w \left\{ \underbrace{\frac{1}{2} \|Xw - y\|^2}_{L(w)} + \underbrace{\alpha \|w\|_1}_{R(w)} \right\}$$

$\text{PROX}_{\tau R}(z) = ?$

$$\left[\left(I + \tau \alpha \delta \| \cdot \|_1 \right)^{-1} (z) \right]_j = \begin{cases} z_j - \tau \alpha, & z_j > \tau \alpha \\ 0, & |z_j| \leq \tau \alpha \\ z_j + \tau \alpha, & z_j < -\tau \alpha \end{cases}$$

(Assignment)

$$w_2^{k+1} = \left(I + \tau \alpha \delta \| \cdot \|_1 \right)^{-1} \left[(I - \tau X^T X) w^k + \tau X^T y \right]$$

Section 6 Deep Learning



Unexpected

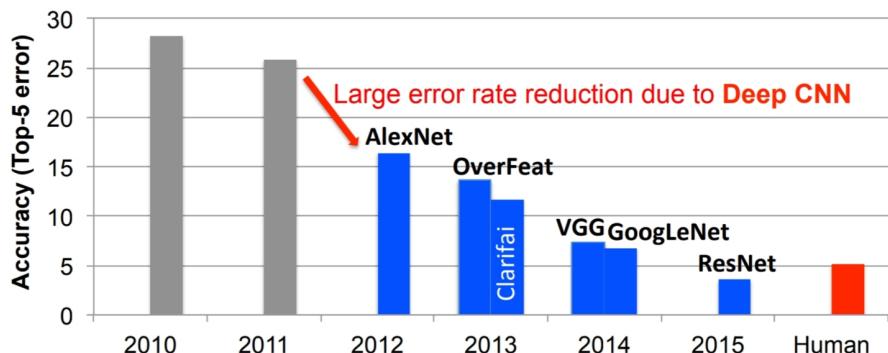


knowledge.net
HowOld.net

Google
amazon.com

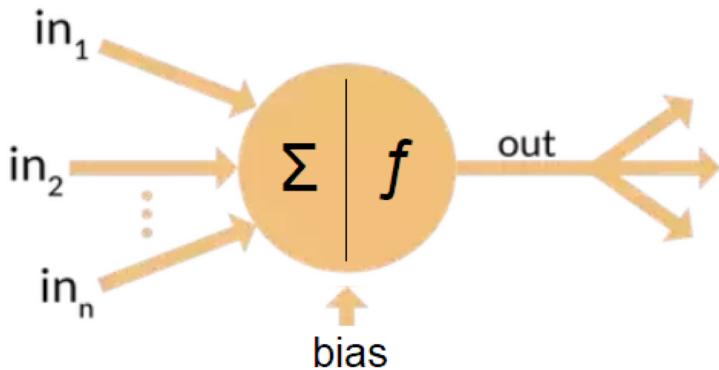
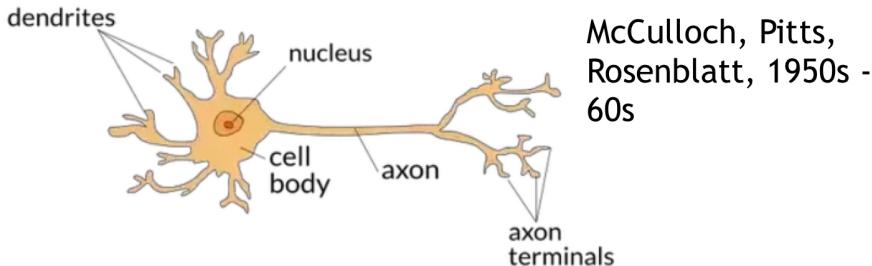


facebook



Speech and image recognition, image captioning, natural language processing

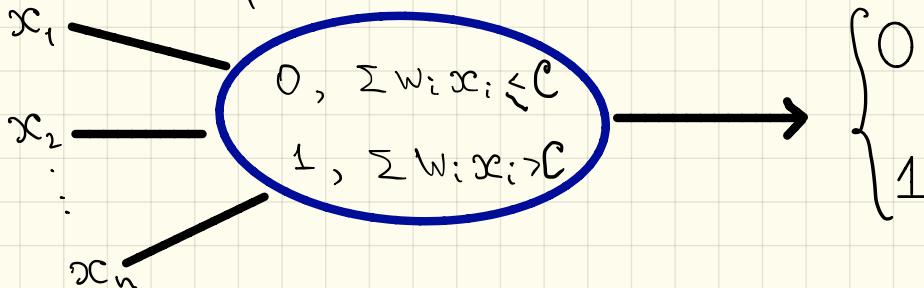
Human neural network



Artificial neural network

Perceptron

Example



x_1 Is the weather good

x_2 Will the partner join

x_3 Is there a public transport connection

$$x_i \in \{0, 1\}$$

$$w = (6, 3, 2)^\top$$

$$C = 5$$

Outcome:

1. The weather is good \rightarrow going
2. The weather is bad \rightarrow not going

In the previous sections we considered simple one step optimisation problems. In deep learning one deals with the objective function being defined as a convolution of several simple functions (called activation functions)

$$f_w(x) = \varphi_L(\varphi_{L-1}(\dots \varphi_1(\varphi_0(x, w_0), w_1), \dots, w_{L-1}), w_L)$$

here $\{\varphi_e\}_{e=1}^L$ are activation functions that are parametrised with weights $w := \{w_e\}_{e=1}^L$.

Def. The above is called a deep neural network with L layers.

What are the typical activation functions?

1. Affine-linear transformations

$$\varphi(x; W, b) = W^T x + b$$

↑ weight matrix ↑ bias vector

$$x \in \mathbb{R}^d, W \in \mathbb{R}^{d \times m}, b \in \mathbb{R}^m$$

2. Heaviside function

$$H(x) = \begin{cases} 1 & , x \geq 0 \\ 0 & , x < 0 \end{cases}, \quad \varphi(x) = (H(x_1), \dots, H(x_d))^T$$

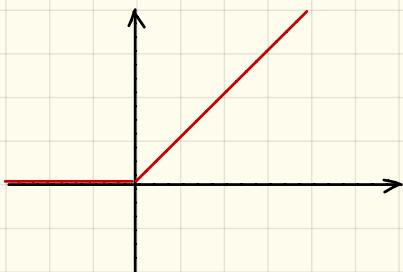
3. Perceptron

$$\varphi(x; W, b) = H(W^T x + b)$$

↑ a two layer neural network itself

4. Rectifier

$$\varphi(x) = (\max(0, x_1), \max(0, x_2), \dots, \max(0, x_d))^T$$



$$\max(0, x) = \text{prox}_{\chi_{[0, +\infty)}}(x)$$

Remark: all the above functions form Rectified Linear Unit (ReLU), and the general function from this set takes a form

$$\varphi(x; W, b) = \max \left\{ 0, \underbrace{\underbrace{W^T x + b}_\text{elementwise}}_\text{vectors in } \mathbb{R}^m \right\}.$$

Typical ReLU neural network is:

$$\varphi(x; w, b) = \max(0, w_1^T \cdot (\max(0, w_2^T \cdots \cdot \max(0, w_l^T x + b_l) \cdots + b_{l-1}) + b_l))$$

$$\varphi(x; w, b) = \max(0, w_1^T x^{(1)} + b^{(1)})$$

$$x^{(e)} = \begin{cases} \max(0, w_e^T x^{(e-1)} + b^{(e)}), & 2 \leq e \leq l \\ \max(0, w_1^T x + b^{(1)}), & e=1 \end{cases}$$

5. Softmax

$$\varphi(x_1, x_2, \dots, x_n) = \left(\frac{e^{x_1}}{\sum e^{x_i}}, \frac{e^{x_2}}{\sum e^{x_i}}, \dots, \frac{e^{x_n}}{\sum e^{x_i}} \right)$$

This can be thought as mapping vectors to probability distributions, indeed for

$$p_k = \frac{e^{x_k}}{\sum e^{x_i}}$$

one has $\sum_{k=1}^n p_k = 1$ and $p_k > 0, k = 1, \dots, n$.

Optimisation problem

$$\hat{w} = \arg \min_w \left\{ \frac{1}{S} \sum_{i=1}^S l(f_w(x_i), y_i) \right\}$$

Example: MSE

$$\hat{w} = \arg \min_w \frac{1}{S} \sum_{i=1}^S \|x_i^{(l)} - y_i\|^2, \text{ where}$$

$$\begin{cases} x_i^{(l)} = \varphi_e(x_i^{(l-1)}, w^{(l)}), & l = 1, \dots, L \\ x_i^{(0)} = x_i \end{cases}$$

Training of deep learning models

Main problem: objective function is not convex.
But we will still apply our tools with a hope for convergence. We can use the gradient descent method in the form of

$$L(W_1, \dots, W_L, b_1, \dots, b_L) = \frac{1}{S} \sum_{i=1}^S l(x_i^{(l)}, y_i)$$

where

$$x_i^{(l)} = \sigma(z_i^{(l)})$$

$$z_i^{(l)} = W_l^\top x_i^{(l-1)} + b_l, \quad l = 1, \dots, L$$

$$x_i^{(0)} = x_i$$

We want to apply the gradient descent and thus need to find the gradient of L with respect to all parameters $W_{jk}^{(l)}, b_j^{(l)}$

Example:

$$L = 2$$

$$d = m_1 = m_2 = 1$$

$$L = \frac{1}{2S} \sum_{i=1}^S \left| \sigma \left(W^{(2)} \left[\sigma \left(W^{(1)} x_i^{(1)} + b^{(1)} \right) \right] + b^{(2)} \right) - y_i \right|^2$$

$$S = 1$$

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial x^{(2)}} \cdot \frac{\partial x^{(2)}}{\partial W^{(2)}} = \frac{\partial L}{\partial x^{(2)}} \cdot x^{(1)}$$

$$\frac{\partial L}{\partial x^{(2)}} = (x^{(2)} - y)$$

$$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial x^{(1)}} \cdot \frac{\partial x^{(1)}}{\partial W^{(1)}} = \frac{\partial L}{\partial x^{(1)}} \cdot x; \quad \frac{\partial L}{\partial x^{(1)}} = \frac{\partial L}{\partial x^{(2)}} \cdot \frac{\partial x^{(2)}}{\partial x^{(1)}} \\ = \frac{\partial L}{\partial x^{(2)}} \cdot \sigma'(z^{(1)}) \cdot W^{(2)}$$

Theorem

Let $f_w(x)$ be a neural network defined as: $f_w(x) = x^{(L)}$

$$x^{(e)} = \begin{cases} x & l=0 \\ \sigma(z^{(e)}), & l=1, \dots, L \end{cases} \quad z^{(e)} = W_e^T x^{(e-1)} + b^{(e)}, \quad l=1, \dots, L$$

Let also the empirical risk function takes the form

$$L(W_1, \dots, W_L, b^{(1)}, \dots, b^{(L)}) = \frac{1}{S} \sum_{i=1}^S \ell(f_w(x_i), y_i)$$

Then the gradient of the empirical risk function with respect to parameters

can be evaluated by a backpropagation scheme:

$$\delta_j^{(L)} := \frac{\partial L}{\partial x_j^{(L)}} = \frac{1}{S} \ell'_x(x^{(L)}, y_i)$$

$$\begin{aligned} \delta_j^{(e)} &:= \frac{\partial L}{\partial x_j^{(e)}} = \sum_{k=1}^d \frac{\partial L}{\partial x_k^{(e+1)}} \cdot \sigma'(z^{(e)}) \cdot (W_{e+1})_{jk} \\ &= \sigma'(z^{(e)}) \cdot [W_{e+1} \delta_j^{(e+1)}] \end{aligned}$$

$$\frac{\partial L}{\partial b_j^{(e)}} = \delta_j^{(e)}$$

$$\frac{\partial L}{\partial w_{jk}^{(e)}} = \sum_j \delta_j^{(e)} x_k^{(e-1)}$$

for $s=1$