

Lecture 5:

Model selection

Gradient descent

The LASSO

The LASSO and gradient descent

Lecture 5

- 1 Model selection: grid search
- 2 Gradient descent method
- 3 The LASSO
- 4 Gradient descent and the LASSO
- 5 Proximal gradient descent
- 6 Deep learning next week

Section 1: Model selection

Main question of the section is: how to select an optimal value of α . The same type of problem arises in polynomial regression as well: how do we select an optimal degree d of a polynomial?

Idea: choose such a value of a hyperparameter that minimizes the error (cost).

Let \mathcal{D} be a probability distribution on $\mathbb{R}^d \times \mathbb{R}^m$ and the input data is built from independent copies $(\bar{x}^{(i)}, \bar{y}^{(i)})$ chosen from the distribution \mathcal{D} .

$$\mathbb{E}(f) := \mathbb{E}_{\bar{x}, \bar{y}} [\ell(\bar{y}, f(\bar{x}))] = \int_{\mathbb{R}^d \times \mathbb{R}^m} p(\bar{x}, \bar{y}) \ell(\bar{y}, f(\bar{x})) d\bar{x} d\bar{y}$$

↑ Population risk
Expected risk
Expected error

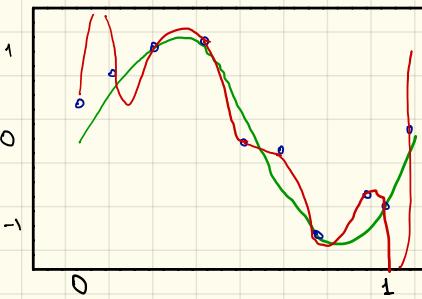
↑ cost function (MSE, MAE, etc.)

$$L_s(f) := \frac{1}{|S|} \sum_{(\bar{x}, \bar{y}) \in S} \ell(\bar{y}, f(\bar{x})) \xrightarrow{|S| \rightarrow \infty} \mathbb{E}(f)$$

↑ Empirical risk

Main problem: f is a function of a training set and thus training error (the one computed over a training set) is not representative.

Example:



The error is zero in this case, but the regression is not optimal.

Thus we need to validate/check the model on a different data.

Idea: split the training data into proper training data and a validation data. Usually 80/20.

$$S = S_t \cup S_v$$

\uparrow training \uparrow validation

$$f = f_{w_t}, \quad w_t = \arg \min_w L_t(f) \quad \text{where}$$

$$L_t(f) := \frac{1}{|S_t|} \sum_{\bar{x}, \bar{y} \in S_t} l(\bar{y}, f(\bar{x}))$$

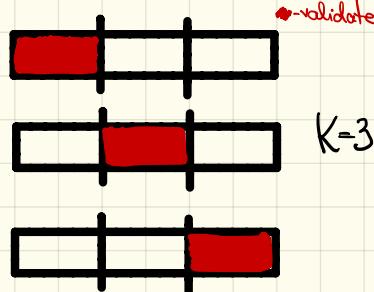
$$\text{Error : } L_v(f) := \frac{1}{|S_v|} \sum_{\bar{x}, \bar{y} \in S_v} l(\bar{y}, f(\bar{x})) \rightarrow \min$$

hyperparameters

K-fold cross-validation

Another possible validation idea: repeat the above several times in such a way that all data is used for both: training and validation but at different times.

- Randomly partition data into K groups
- Train the model K times, leaving 1 group for validation
- Average the result



1. Split the data S into training data S_t and validation data S_v .
2. Train the model.
3. Calculate the validation error.
4. Minimise the validation error over all possible values of hyperparameters.

Example 1: Ridge regression - selection of a parameter α .

$$W_d = \arg \min_{W \in \mathbb{R}^{(d+1) \times m}} \frac{1}{2} \| X_t W - Y_t \|^2 + \frac{\alpha}{2} \| W \|^2$$

$$\hat{\alpha} = \arg \min_{\alpha} \frac{1}{2} \| X_v W^d - Y_v \|^2$$

↑
 $X_t \in \mathbb{R}^{(d+1) \times m}$
 ↑
 $Y_t \in \mathbb{R}^{m \times 1}$
 ↑
 $X_v \in \mathbb{R}^{(d+1) \times m}$
 ↑
 $Y_v \in \mathbb{R}^{m \times 1}$

Example 2: Polynomial regression - selection of a parameter d .

$$\hat{W}_d = \arg \min_{W \in \mathbb{R}^{(d+1) \times m}} \frac{1}{2} \| \Phi_d(X_t) W - Y_t \|^2$$

$$\hat{d} = \arg \min_d \frac{1}{2} \| \Phi_d(X_v) W_d - Y_v \|^2$$

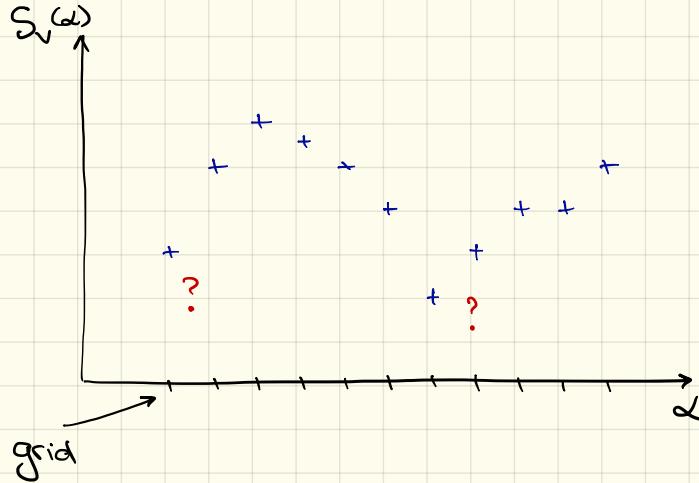
? How can we solve such a bi-level optimisation problem?

1. Analytically it may be too complicated.
2. Grid search (evaluate the function of the interest - validation error - at the points of a grid).
3. Random search (evaluate the function of our interest at a random points).
4. Gradient-based optimisation (iterative procedure to approach the minimum)
5. Bayesian optimisation
6. Evolutionary optimisation

In this module we are mainly using grid search.

The method consist of:

- defining a grid of values of hyper parameter;
- evaluating the validation error for every value of a hyper parameter of the grid;
- finding such a λ that minimises the error.



- Where the minimum is?
- Whether it is attained at one of the grid points?
- Can we be sure that this is a global minimum and not the local one?

Advantages:

1. Universality
2. Easy realisation

Disadvantages

1. Computationally infeasible
2. No guarantee of getting the minimum point.

Section 2: The LASSO

Previously we have seen how the ridge regression can help to prevent overfitting.

$$W_\alpha = \arg \min_w \left\{ L(W) + R(W) \right\}$$

↑ regression ↑ regularisation

In this section we discuss another regularisation of a standard regression problem. Namely we discuss the LASSO - Least Absolute Shrinkage and Selection Operator.

Def. The LASSO problem is a minimisation problem of the form

$$W_\alpha = \arg \min_w \frac{1}{2} \| \Phi(X)W - Y \|^2 + \alpha \| W \|_1$$

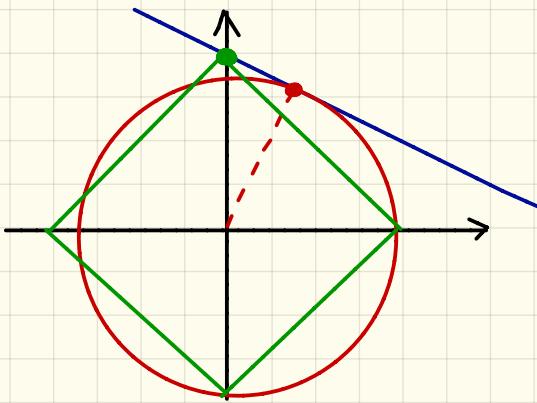
$$\text{where } \| W \|_1 = \sum_{i,j} |W_j^{(i)}|$$

What is the motivation of introducing 1-norm?

$$\| \boldsymbol{\sigma} \|_2^2 = \sum_{i=1}^n \sigma_i^2 \rightarrow \text{Fix} \Rightarrow \| \boldsymbol{\sigma} \|_1 \rightarrow \min \text{ when } \boldsymbol{\sigma} \text{ is sparse}$$

Example:

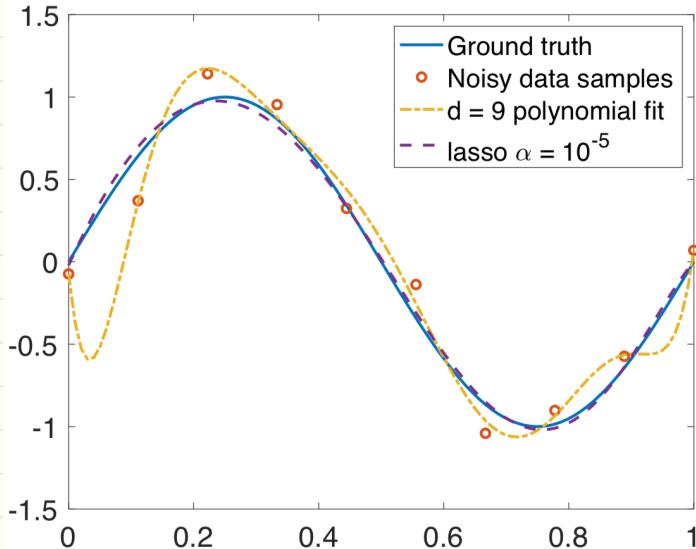
Consider a linear regression with 1 data point:
 (x, y) .



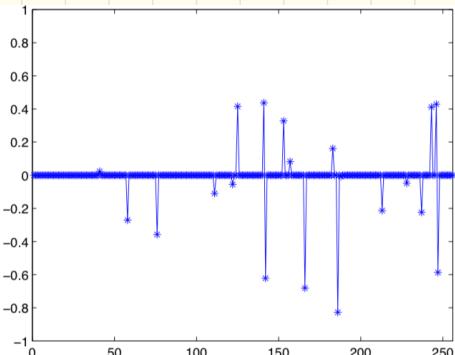
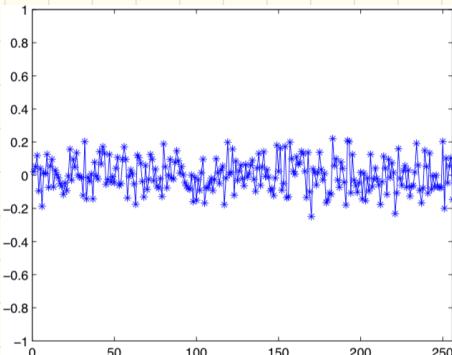
$$\stackrel{1^{\circ}}{=} \min(w_0^2 + w_1^2) \\ \text{subject to} \\ y = w_0 + w_1 x$$

$$\stackrel{2^{\circ}}{=} \min(|w_0| + |w_1|) \\ \text{subject to} \\ y = w_0 + w_1 x$$

Example:



Example:



$$\|\cdot\|_1 = 20.061$$

$$\|\cdot\|_1 = 6.293$$

$$\|\cdot\|_2 = 1.5431$$

$$\|\cdot\|_2 = 1.747$$

Therefore, the LASSO problem reconstructs solutions with only few non-zero elements. This leads to implicit parameters reduction and thus we obtain the simpler model (occam's razor).

Problem: 1-norm is not differentiable.

We need another method to solve the problem!

Section 3: Gradient descent

Let us consider a problem of minimisation

$$\hat{w} = \arg \min_w E(w)$$

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

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

for the energy E , an initial value w^0 and a step-size parameter τ .

$$w^1 = w^0 - \tau \nabla E(w^0)$$

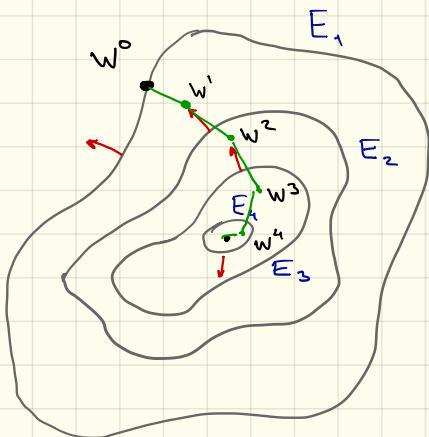
$$w^2 = w^1 - \tau \nabla E(w^1)$$

$$w^3 = w^2 - \tau \nabla E(w^2)$$

.

.

.



$$E_1 > E_2 > E_3 > E_4$$

Example Zero order polynomial regression

$$E(w) = \frac{1}{2s} \sum_{i=1}^s (w - y_i)^2 \Rightarrow \nabla E(w) = w - \frac{1}{s} \sum_{i=1}^s y_i$$

$$w^{k+1} = (1-\tau) w^k + \frac{1}{s} \sum_{i=1}^s y_i$$

what happens for $\tau = 1$?

$$(w^{k+1} - \frac{1}{s} \sum_{i=1}^s y_i) = (1-\tau) (w^k - \frac{1}{s} \sum_{i=1}^s y_i)$$
$$0 < \tau < 2 \Rightarrow w^k \rightarrow \frac{1}{s} \sum_{i=1}^s y_i$$

Example Linear regression

$$E(w) = \frac{1}{2s} \|Xw - \bar{y}\|^2$$

$$\nabla E(w) = \frac{1}{s} X^T (Xw - \bar{y})$$

$$\begin{aligned} w^{k+1} &= w^k - \frac{\tau}{s} X^T (Xw^k - \bar{y}) \\ &= \left(1 - \frac{\tau}{s} X^T X\right) w^k + \frac{\tau}{s} X^T \bar{y} \end{aligned}$$

How should we select the parameter τ to guarantee the convergence?

Idea: we should guarantee that the energy decrease!

$$\begin{aligned} E(w^{k+1}) &= \frac{1}{2s} \|Xw^k - \frac{\tau}{s} X^T X w^k + \frac{\tau}{s} X^T \bar{y} - \bar{y}\|^2 \\ &= \frac{1}{2s} \left\| \left(I - \frac{\tau}{s} X^T X\right) \cdot (Xw^k - \bar{y}) \right\|^2 \\ &\leq \frac{1}{2s} \left\| I - \frac{\tau}{s} X^T X \right\|^2 \cdot \|Xw^k - \bar{y}\|^2 \\ &= \left\| I - \frac{\tau}{s} X^T X \right\|^2 \cdot E(w^k) \end{aligned}$$

↑↑ for small τ

How to guarantee a convergence in general case?

Theorem: Let $E: C \subset \mathbb{R}^m \rightarrow \mathbb{R}$ be a convex energy function such that $G(w) := \frac{1}{2\tau} \|w\|^2 - E(w)$ is also convex. Let \hat{w} be a global minimiser of E . Then:

1. For any k : $E(w^{k+1}) \leq E(w^k)$
2. There exist $C > 0$ such that $E(w^k) - E(\hat{w}) \leq C/k$
3. As a consequence $E(w^k) \rightarrow E(\hat{w})$

Def.: A continuously differentiable function $E: C \subset \mathbb{R}^m \rightarrow \mathbb{R}$ is called L-smooth if

$$\|\nabla E(u) - \nabla E(v)\| \leq L \|u - v\|, \quad \forall u, v \in C.$$

Lemma. Let $E: \mathbb{C} \subset \mathbb{R}^m \rightarrow \mathbb{R}$ be an γ -smooth function.

Then

$$G_\tau(w) = \frac{1}{2\tau} \|w\|^2 - E(w)$$

is convex on \mathbb{C} .

Example

$$E(w) = \frac{1}{2s} \sum_{j=1}^s (w - y_j)^2$$

$$\nabla E(w) = w - \frac{1}{s} \sum_{j=1}^s y_j$$

$$\|\nabla E(u) - \nabla E(v)\| = \|u - v\|$$

$\Rightarrow E$ is 1 - smooth ($\tau = 1$)

What if E is not convex?

Then the gradient descent would end up at some local minima.

What if E is convex but there is no such a $\tilde{\tau}$ that G_τ is convex?

Then we need to adapt our algorithm. When defining w^{k+1} one needs to check $E(w^{k+1}) \leq E(w^k)$. If not, then decrease the value of τ .

Section 4: The LASSO and gradient descent

How to use a gradient descent to solve the LASSO?

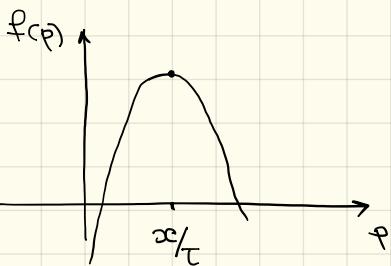
The main difficulty to overcome is non-differentiability of the absolute value function. The trick is to replace the absolute value with a differentiable function.

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

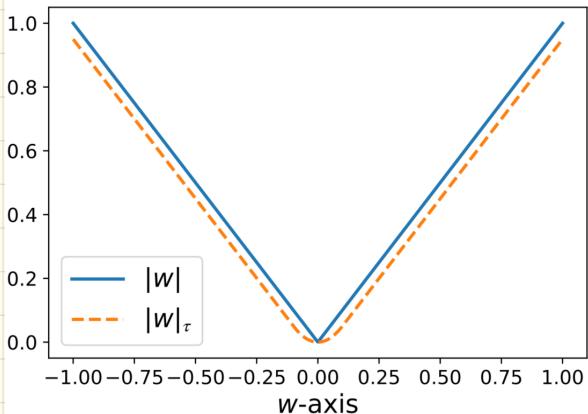
$$|x|_\tau := \max_{p \in [-1, 1]} xp - \frac{\tau}{2} p^2, \quad \tau > 0$$

↑ small perturbation

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



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



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

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

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

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

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

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

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

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

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

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

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

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

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

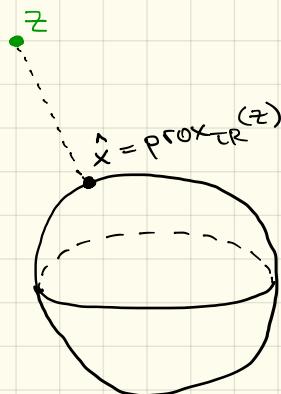
$$\text{Example: } R(w) = \begin{cases} 0, & w \in C \\ \infty, & w \notin C \end{cases} \quad \text{↑ convex set}$$

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

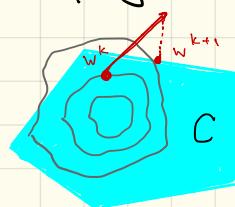
$$\text{prox}_{\tau R}(z) = \arg \min_x \left[\frac{1}{2} \|x - z\|^2 + \tau R(x) \right]$$

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

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)

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