# 24/12/16 Lasso / Ridge regularization

$L^p$ norm (=P-norm): Method which is used to gauge magnitude or distance of vector in $LP$ space (=Lebesgue space)

• $L^p$ space (=Lebesgue space): function space where is used to gauge the size of function.

Especially, this space is defined as the set of functions whose absolute value raised to the power of $p$ (=$f^p$) is integrable. (=함의 크기를 정의)

↳ $L^p(\Omega) = \left\{ f : \int_\Omega |f(x)|^p dx < \text{inf.} \right\}$, s.t $1 \leq p < \infty$

if $p=1 \Rightarrow$ use $L^1$ norm
if $p=2 \Rightarrow$ use $L^2$ norm ...

$\Rightarrow L^p$ norm (=P norm) $= \|X\|_p = \left( \sum_i |X_i|^p \right)^{1/p}$, s.t $1 \leq p < \infty$ (=함의 크기를 측정)

tool.

---

Lasso (=L1): Least Absolute Shrinkage and Selection Operator $\Rightarrow$ Passive operator

↳ Lasso is seemed select a few features which are relatively important.

↳ Lasso tries to minimize absolute value of weight value (=$opt(\sum_i |W_i|)$)

loss function in Lasso regularization = MSE + $\lambda|W| = L_1$

$= \frac{1}{n} \sum_{i=1}^{n} (\hat{Y}_i - Y_i)^2 + \frac{\lambda}{2n} \sum_{i=1}^{n} |W_i|$,

$\frac{\partial L_1}{\partial w} = \frac{\partial}{\partial w}(Wx+b-y)^2 + \lambda \frac{|W|}{\partial w} = \begin{cases} 2x(Wx+b-y) + \lambda & \text{s.t. } w \geq 0 \\ 2x(Wx+b-y) - \lambda & \text{s.t. } w < 0 \end{cases}$

∴ let $A = 2x(Wx+b-y)$, then $W^* = \begin{cases} W + y \cdot A - \lambda & \text{s.t. } w \geq 0 \\ W + y \cdot A + \lambda & \text{s.t. } w < 0 \end{cases}$

$\Rightarrow$ Closer to Zero, and some values are even equal to zero.

↳ Effect of regulation

∵ every value are subtracted by $\lambda$ equally.
(= $\lambda$ 계틍 배뼝 연개간 뫛일 W퓼는 0이 됨)

$\Rightarrow$ If assuming that all weight values have same scale, small weight values which are less important will be deleted. (정라 강둥 낭고 개깅연4 상영서닌 없어나값)

---

Ridge (=L2) loss function = MSE + $\lambda W^2 = L_2 \Rightarrow$ Active operator

$= \frac{1}{n} \sum_{i=1}^{n} (\hat{Y}_i - Y_i)^2 + \frac{\lambda}{2n} \sum_{i=1}^{n} W_i^2$,

$\frac{\partial L_2}{\partial w} = \frac{\partial}{\partial w}(Wx+b-y)^2 + \lambda \frac{W^2}{\partial w} = 2x(Wx+b-y) + 2\lambda W$

→ Effect of regulation

∴ let $A = 2x(Wx+b-y)$, then $W^* = W - y \cdot (A + 2\lambda W)$ $\Rightarrow$ It operates similarly with $L_1$, but it regularize $W^*$ more rationally.
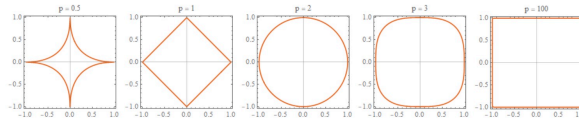
if $W$ is big relatively, $L_2$ will subtract with bigger value, but if $W$ is small, $L_2$ will subtract with smaller value. $\Rightarrow$ Sensitively Activating to Outlier than $L_1$

cf) root is eliminated for enhancing efficiency of calculations.

Dif.

$\otimes$ Distribution of P-norm in 2-dimension vector space



$\otimes$ Distribution of P-norm in 3-dimension vector space



Ref.

# +) Why do we want to make weight values close to zero?

= (How large weights are related to overfitting)

If weights are small, the features have a small effect on prediction. the model fits the general trend of the data. Else if weights are large, the model becomes highly sensitive to small changes to small changes in specific features.

In linear regression model, the prediction is:

$\hat{y} = \sum W_i x_i + b$, If $W_i$ has large value, a small change in $x_i$ will cause a large change in $\hat{y}$.