# Normal Equation and Least Squares

The least Squares problem:

$$\min_{\beta} \ell(\beta) = \| y - X\beta \|_2^2$$

$$= y^T y - 2 y^T X \beta + \beta^T X^T X \beta$$

$X \in \mathbb{R}^{n \times p}$ data matrix

$y \in \mathbb{R}^{n \times 1}$ target values

$n$ : number of sample points

$p$ : dimension of feature vectors.

Solve by setting the gradient to zero:

$$\nabla_\beta \ell(\beta) = -2 X^T y + 2 X^T X \beta = 0$$

$$\Leftrightarrow \quad X^T X \beta = X^T y \quad , \quad \text{called the "normal equation."}$$

If $X^T X$ is invertible, $\hat\beta = (X^T X)^{-1} X^T y$ is the unique solution.

Q: Under what condition is $(X^T X)$ invertible, or equivalent, of full rank?

Note: The rank of a square matrix is the max # of linearly independent rows (or columns).

A: Two cases: ① $n < p$ ② $n \geq p$.



$n < p$

$n \geq p$

$\text{rank}(X^T X) \leq n$,

because every column of $X^T X$ $\Rightarrow$ is a linear combination of at most $n$ $p$-dimensional vectors.

When $n \leq p$, rank$(X^TX) < p$, so $X^TX$ not invertible,
and the least square problem has multiple solutions.

When $n \geq p$, and there are $p$ linearly independent feature vectors in the
data, (which is usually the case when $n > p$), $X^TX$ is invertible and

$$\hat{\beta} = (X^TX)^{-1}X^Ty$$ is the unique solution.

# Ridge Regression

$$\min_{\beta} \ell_{ridge}(\beta) = \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

$$= y^Ty - 2y^TX\beta + \beta^T (X^TX + \lambda I) \beta$$

$\lambda > 0$ : regularization parameter,

$I$ : $p$-by-$p$ identity matrix

Solve $\nabla_{\beta} \ell_{ridge}(\beta) = 0$

$$\Leftrightarrow \quad -2X^Ty + 2(X^TX + \lambda I)\beta = 0$$

$$\Leftrightarrow \quad (X^TX + \lambda I)\beta = X^Ty.$$

Thm: $X^TX + \lambda I$ is always invertible

pf: Prove the following lemma first:

> Lemma: $\forall a \in \mathbb{R}^p$, $a$ not the zero vector,
> $$a^T (X^TX + \lambda I) a > 0.$$
> pf: $a^T(X^TX + \lambda I) a = a^TX^TXa + \lambda a^Ta$
> $$= \|Xa\|_2^2 + \lambda a^Ta \quad > 0. \text{ since } a \neq 0 \text{ and } \lambda > 0$$

Then prove by contradiction: If $X^TX + \lambda I$ is not invertible, its columns are not $\underset{\text{linearly}}{}$ independent, so there exists $a \in \mathbb{R}^p$, $a \neq 0$ such that

$$(X^TX + \lambda I)a = 0,$$

which implies $a^T(X^TX + \lambda I)a = 0$, a contradiction to the lemma.

#

So

$$\hat{\beta}_{ridge} = (X^T X + \lambda I)^T X^T y$$

is the unique solution to the Ridge Regression problem.

Why ridge regression?

① When $n < p$, helps to get a unique solution.

② When $n \geq p$, even though $\hat{\beta}$ usually exists and is unique, it may overfit the data. In terms of Bias and variance,

$$bias(\hat{\beta}_{ridge}) \geq bias(\hat{\beta}) = 0 \text{ under the linear model,}$$

$$Variance(\hat{\beta}_{ridge}) < Variance(\hat{\beta})$$

As $\lambda \uparrow$, $bias(\hat{\beta}_{ridge}) \uparrow$ and $Variance(\hat{\beta}_{ridge}) \downarrow$

Use cross validation to decide $\lambda$.

## Histogram.

Consider the following family of p.d.f.s over the 1-d interval $[a,b]$:

$$f(x) = \sum_{j=1}^{k} \mathbb{1}\{x \in Bin_j\} \, P_j \quad , \quad P_j \geq 0 \text{ is the density in the } j^{th} \text{ bin.}$$

Let $\Delta_1, \Delta_2, \dots, \Delta_k$ be the $\underset{\text{pre-specified}}{\text{sizes}}$ of the $k$ bins, so $\sum_{j=1}^{k} \Delta_j = b-a$

and $\quad \text{Prob}(X \in Bin_j) = \int_a^b \mathbb{1}\{x \in Bin_j\} f(x) \, dx = P_j \Delta_j.$

Since $f(x)$ is a p.d.f, we have

$$\int_a^b f(x) \, dx = \sum_{j=1}^{k} P_j \Delta_j = 1$$

Given an i.i.d sample $\{x_1, x_2, \dots x_n\}$ drawn from some $f$ in this family, we want to estimate the densities $P_1, P_2, \dots, P_k$. We do ML estimation.

Likelihood : $\quad L(P_1 \dots P_k) = \prod_{i=1}^{n} \prod_{j=1}^{k} (P_j \Delta_j)^{\mathbb{1}\{x_i \in Bin_j\}}$

Log likelihood: $\quad \ell(P_1 \dots P_k) = \sum_{i=1}^{n} \sum_{j=1}^{k} \mathbb{1}\{x_i \in Bin_j\} \log(P_j \Delta_j)$

$$= \sum_{j=1}^{k} \underbrace{\sum_{i=1}^{n} \mathbb{1}\{x_i \in Bin_j\}}_{} \log(P_j \Delta_j)$$

call $n_j$, # of points in $Bin_j$

$$\text{concave in } P_1 \dots P_k$$

Solve $\quad \max \ell(P_1, \dots P_k) \quad \text{s.t.} \quad \sum_j P_j \Delta_j = 1 \quad$ by setting the gradient

of the Lagrangian function to zero :

$$\text{\circled{1}} \; \partial_{P_{j'}} \left[ \ell(P_1 \dots P_k) - \lambda \left( \sum_j P_j \Delta_j - 1 \right) \right] = 0 \quad \Longleftrightarrow \quad \frac{n_{j'}}{P_{j'}} - \lambda \Delta_{j'} = 0$$

$$\therefore P_{j'} = \frac{n_{j'}}{\lambda \Delta_{j'}}. \qquad \text{Since } \sum_{j'} P_{j'} \Delta_{j'} = 1, \quad \lambda \text{ must be } \sum_{j'} n_{j'} = n \text{, and}$$

$$P_{j'} = \frac{n_{j'}}{n \Delta_{j'}}, \quad \text{the histogram density estimate.}$$
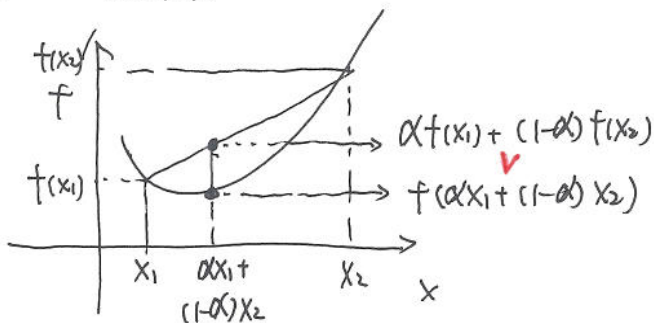
# $L_0$ penalty is non-convex

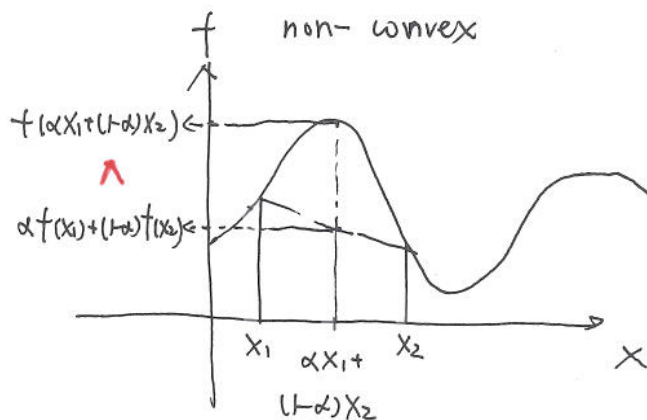For simplicity, consider one dimensional case.

Def. A function $f$ is convex if

$$\alpha f(x_1) + (1-\alpha) f(x_2) \geq f(\alpha x_1 + (1-\alpha) x_2) \qquad \forall \; 0 \leq \alpha \leq 1 \quad \text{and} \quad \forall \; x_1, x_2 \text{ in domain of } f.$$

Ex.  convex



The $L_0$ penalty in 1-d:

$$L_0(\beta) = \mathbb{1}\{\beta \neq 0\}$$



$$f(\alpha x_1 + (1-\alpha) x_2) = 1$$
$$\alpha f(x_1) + (1-\alpha) f(x_2) = 1 - \alpha < 1.$$

$\Rightarrow$ non-convex.