

Element of Statistical Learning Note

Zichong Wang

January 9, 2026

Contents

1	Intro	3
2	Chap 3: Linear Methods for Regression	4
2.1	Confidence region for $\hat{\beta}$	4
2.2	What is Linear	4
2.3	Gauss-Markov Theorem	4
2.4	QR decomposition	5
2.4.1	Application to Least Squared	5
2.4.2	About orthogonal matrix	6
2.4.3	SVD and orthogonal matrix	7
2.5	Multiple testing in forward selection	8
2.6	Ridge regression	8
2.6.1	How to use ridge regression	8
2.6.2	Equivalence of two forms	8
2.6.3	Bayesian view	10
2.6.4	SVD in ridge	10
2.6.5	Revisit PCA	11
2.7	LASSO	12
2.7.1	Subgradient and subdifferential	12
2.7.2	Why shrink to zero	13
2.7.3	Soft thresholding	14
3	Chap 4: Linear methods for classification	15
3.1	Linear Regression of an Indicator Matrix	15

4	Optimization methods	17
4.1	Norms of matrices	17
4.1.1	Induced norm	17
4.1.2	Frobenius norm	18
4.1.3	Schatten norm	19
4.2	Gradient descent	19
4.2.1	Descent Lemma	19
4.2.2	Convergence rate	20
4.2.3	Another view of GD	20
4.3	Subgradient descent	21
4.4	Proximal Gradient Descent	21
4.4.1	ISTA	22
4.5	Stochastic gradient descent	22

1 Intro

While reading ESL, I felt the need to write things down—to record what I did not understand and how I eventually worked through it. This note grew out of that process. It can be regarded as a personal supplementary material of ESL. It covers concepts from analysis, linear algebra, optimization, and computation, with a primary emphasis on statistics.

2 Chap 3: Linear Methods for Regression

2.1 Confidence region for $\hat{\beta}$

We know that $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$, thus $\hat{\beta} - \beta \sim \mathcal{N}(0, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$.

Since for $\mathbf{z} \sim \mathcal{N}(0, \mathbf{\Sigma})$, we have $\mathbf{z}^T \mathbf{\Sigma}^{-1} \mathbf{z} \sim \chi_k^2$, where $k = \text{rank}(\mathbf{\Sigma})$, we have

$$(\hat{\beta} - \beta)^T (\sigma^2(\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\hat{\beta} - \beta) \sim \chi_{p+1}^2$$

Thus,

$$\frac{1}{\sigma^2} (\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta) \sim \chi_{p+1}^2$$

And have approx confidence region

$$(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi_{p+1, 1-\alpha}^2$$

Actually, $\hat{\sigma}^2 = \frac{1}{N-p-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N-p-1} \text{RSS}$, and $\text{RSS}/\sigma^2 \sim \chi_{N-p-1}^2$,

$$\frac{(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta)/(p+1)}{\hat{\sigma}^2} \sim F_{p+1, N-p-1}$$

2.2 What is Linear

In the context of **linear model**, we are talking about linearity in parameters, meaning that the prediction \hat{y} is a linear combination of the parameters β_j . The \mathbf{X} itself can be non-linear transformations of the original features, e.g., polynomial terms, interaction terms, etc. $y = 1/(\beta_0 + \beta_1 x)$ and $y = \beta_0 e^{\beta_1 x}$ are not linear models, since they're not linear in parameters.

In the context of **Linear estimators**, we are talking about the estimator $\hat{\theta}$ (e.g. $\hat{\beta}$) can be written as a linear combination of the observed response values y_i , i.e. $\hat{\theta} = \mathbf{c}^T \mathbf{y}$. The weight \mathbf{c} depends only on \mathbf{X} , not on \mathbf{y} . A linear estimator **can be** a prediction at a new point, or the estimated coefficients $\hat{\beta}$ themselves.

2.3 Gauss-Markov Theorem

Why assume only know \mathbf{X} , but not \mathbf{y} ?

Note that though y_i as sample responses, are observable, the following statements and arguments including assumptions, proofs and the others assume under the only condition of knowing $\mathbf{X}_{i,j}$ but not y_i . — [2]

We have a *challenger* linear estimator $\tilde{\beta} = \mathbf{C} \mathbf{y}$, where $\mathbf{C} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D}$, a modification

of OLS estimator. Ensure it's unbiased:

$$\begin{aligned}
\mathbb{E}(\tilde{\beta}) &= \mathbb{E}(\mathbf{C}\mathbf{y}) \\
&= \mathbf{C}\mathbb{E}(\mathbf{y}) = ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D})\mathbf{X}\beta \\
&= \beta + \mathbf{D}\mathbf{X}\beta \\
&= \beta
\end{aligned}$$

Meaning $\mathbf{D}\mathbf{X} = 0$.

Now, compute the variance:

$$\begin{aligned}
\text{Var}(\tilde{\beta}) &= \text{Var}(\mathbf{C}\mathbf{y}) \\
&= \mathbf{C}\text{Var}(\mathbf{y})\mathbf{C}^T \\
&= \sigma^2 \mathbf{C}\mathbf{C}^T \\
&= \sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D})((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D})^T \\
&= \sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1} + \mathbf{D}\mathbf{D}^T) \\
&= \text{Var}(\hat{\beta}) + \sigma^2 \mathbf{D}\mathbf{D}^T \\
&\geq \text{Var}(\hat{\beta})
\end{aligned}$$

2.4 QR decomposition

Any real squared matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is orthogonal matrix ($\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$), and \mathbf{R} is upper triangular matrix.

Any rectangular matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m \geq n$), we can decompose it as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is $m \times m$ orthogonal matrix, and \mathbf{R} is $m \times n$ upper triangular matrix (the last $m - n$ rows are all zero). It can be regarded as $\mathbf{A} = \mathbf{Q}\mathbf{R} = [\mathbf{Q}_1 \quad \mathbf{Q}_2] \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \mathbf{Q}_1 \mathbf{R}_1$, where $\mathbf{Q}_1 \in \mathbb{R}^{m \times n}$, $\mathbf{Q}_2 \in \mathbb{R}^{m \times (m-n)}$, $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$ which is upper triangular.

If \mathbf{A} have k linearly independent columns, then first k columns of \mathbf{Q} form an orthonormal basis of the column space of \mathbf{A} . The fact that any column k of \mathbf{A} only depends on the first k columns of \mathbf{Q} corresponds to the triangular form of \mathbf{R} .

QR decomposition can be calculated using Gram-Schmidt process, or using Householder reflections. In practice, Householder reflections are more stable and efficient.

2.4.1 Application to Least Squared

In linear least squares problems, we aim to find a vector \mathbf{x} that minimizes the Euclidean norm of the residual for an overdetermined system $\mathbf{A}\mathbf{x} \approx \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $m \geq n$. The goal is to solve:

$$\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$$

Using the Full QR Decomposition, we substitute $\mathbf{A} = \mathbf{Q} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$:

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 = \left\| \mathbf{Q} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} \mathbf{x} - \mathbf{b} \right\|_2$$

Since multiplying by an orthogonal matrix Q preserves the Euclidean norm, we can left-multiply the entire expression by Q^T :

$$\|Ax - b\|_2 = \left\| \begin{bmatrix} R_1 \\ 0 \end{bmatrix} x - Q^T b \right\|_2$$

If we partition $Q^T b$ into two components— $c_1 \in \mathbb{R}^n$ and $c_2 \in \mathbb{R}^{m-n}$:

$$\left\| \begin{bmatrix} R_1 x - c_1 \\ -c_2 \end{bmatrix} \right\|_2^2 = \|R_1 x - c_1\|_2^2 + \|c_2\|_2^2$$

To minimize the total error, we must make the first term zero. The least squares solution is found by solving the square, upper-triangular system:

$$R_1 x = c_1$$

Since R_1 is upper triangular, we can efficiently solve this system using back substitution, **that's how we solve linear equations manually in algebra class!** The remaining term $\|c_2\|_2$ represents the minimum residual norm (the "error" of the fit).

2.4.2 About orthogonal matrix

Let $Q \in \mathbb{R}^{n \times n}$ be an orthogonal matrix. Write Q in terms of its column vectors:

$$Q = [q_1 \ q_2 \ \cdots \ q_n], \quad q_i^\top q_j = \delta_{ij}.$$

Then

$$Q^\top = \begin{bmatrix} q_1^\top \\ q_2^\top \\ \vdots \\ q_n^\top \end{bmatrix}.$$

Computation of $Q^\top Q$.

$$Q^\top Q = \begin{bmatrix} q_1^\top \\ q_2^\top \\ \vdots \\ q_n^\top \end{bmatrix} [q_1 \ q_2 \ \cdots \ q_n] = \begin{bmatrix} q_1^\top q_1 & q_1^\top q_2 & \cdots & q_1^\top q_n \\ q_2^\top q_1 & q_2^\top q_2 & \cdots & q_2^\top q_n \\ \vdots & \vdots & \ddots & \vdots \\ q_n^\top q_1 & q_n^\top q_2 & \cdots & q_n^\top q_n \end{bmatrix}.$$

Using $q_i^\top q_j = \delta_{ij}$,

$$Q^\top Q = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = I.$$

Computation of QQ^\top .

$$QQ^\top = [q_1 \ q_2 \ \cdots \ q_n] \begin{bmatrix} q_1^\top \\ q_2^\top \\ \vdots \\ q_n^\top \end{bmatrix} = q_1 q_1^\top + q_2 q_2^\top + \cdots + q_n q_n^\top = \sum_{k=1}^n q_k q_k^\top.$$

For any $\mathbf{x} \in \mathbb{R}^n$,

$$(\mathbf{Q}\mathbf{Q}^\top)\mathbf{x} = \sum_{k=1}^n \mathbf{q}_k(\mathbf{q}_k^\top \mathbf{x}) = \mathbf{x},$$

since $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ is an orthonormal basis of \mathbb{R}^n . Hence

$$\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}.$$

Therefore, for an orthogonal matrix \mathbf{Q} ,

$$\mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}\mathbf{Q}^\top = \mathbf{I}.$$

2.4.3 SVD and orthogonal matrix

There are two kinds of SVD: full SVD and reduced SVD. Let $\mathbf{X} \in \mathbb{R}^{N \times p}$, and nomally $N > p$.

- **Full SVD** in standard linear algebra.

$$\mathbf{X} = \mathbf{U}_{\text{full}} \mathbf{D}_{\text{full}} \mathbf{V}^\top$$

- $\mathbf{U}_{\text{full}} \in \mathbb{R}^{N \times N}$ is orthogonal matrix, whose columns are eigenvectors of $\mathbf{X}\mathbf{X}^\top$.
- $\mathbf{D}_{\text{full}} \in \mathbb{R}^{N \times p}$ is a rectangular diagonal matrix, and the last $N - p$ rows are all zero.
- $\mathbf{V} \in \mathbb{R}^{p \times p}$ is orthogonal matrix, whose columns are eigenvectors of $\mathbf{X}^\top \mathbf{X}$, and are basis of row space of \mathbf{X} .

- **Reduced SVD** in linear regression. We simply ignore the zero parts of \mathbf{U} and \mathbf{D} .

\mathbf{D} and \mathbf{U} can be seen as

$$\mathbf{D}_{\text{full}} = \left[\begin{array}{ccc|ccc} d_1 & 0 & \dots & & & \\ 0 & \ddots & & & & \\ \dots & & & d_p & & \\ \hline 0 & 0 & 0 & & & \\ \vdots & \vdots & \vdots & & & \\ 0 & 0 & 0 & & & \end{array} \right] \left\{ \begin{array}{l} p \text{ (Non-zero part)} \\ N - p \text{ (All zero part)} \end{array} \right.$$

$$\mathbf{U}_{\text{full}} = [\mathbf{U}_1 \quad | \quad \mathbf{U}_2]$$

Thus,

$$\begin{aligned} \mathbf{U}_{\text{full}} \cdot \mathbf{D}_{\text{full}} &= [\mathbf{U}_1 \quad \mathbf{U}_2] \cdot \begin{bmatrix} \mathbf{D}_p \\ \mathbf{0} \end{bmatrix} \\ &= \mathbf{U}_1 \cdot \mathbf{D}_p + \mathbf{U}_2 \cdot \mathbf{0} \\ &= \mathbf{U}_1 \cdot \mathbf{D}_p \end{aligned}$$

That's the reduced SVD:

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top$$

where $\mathbf{U} \in \mathbb{R}^{N \times p}$ with orthonormal columns, $\mathbf{D} \in \mathbb{R}^{p \times p}$ diagonal with positive entries, and $\mathbf{V} \in \mathbb{R}^{p \times p}$ orthogonal.

In regression, we usually use reduced SVD, and $\mathbf{U}\mathbf{U}^\top = \mathbf{H} \neq \mathbf{I}$, but $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$. For $\mathbf{y} \in \mathbb{R}^N$, $\mathbf{U}^\top \mathbf{y}$ map \mathbf{y} from \mathbb{R}^N to \mathbb{R}^p , which is the coefficients in the basis of columns of \mathbf{U} . And $\mathbf{U}\mathbf{U}^\top \mathbf{y}$ project \mathbf{y} onto the column space of \mathbf{X} .

2.5 Multiple testing in forward selection

ESL page 60:

Other more traditional packages base the selection on F -statistics, adding “significant” terms, and dropping “non-significant” terms. These are out of fashion, since they do not take proper account of the multiple testing issues.

Assume we have p candidate features, and already selected k features. When considering adding a new feature, we are actually performing $p - k$ hypothesis tests (each test for one feature). Even the rest $p - k$ features are all noise, with significance level α , we still have a probability of $1 - (1 - \alpha)^{p-k}$ to incorrectly add at least one noise feature.

2.6 Ridge regression

Answer questions: why two forms are equivalent? Why not equivariant under scaling of the inputs? What is a good practice for it? df of ridge? In the case of orthogonal inputs, why $\hat{\beta}^{\text{ridge}} = \hat{\beta}/(1 + \lambda)$?

Ridge regression shrinks the regression coefficients by imposing a penalty on their size:

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

and is equivalent to

$$\begin{aligned} \hat{\beta}^{\text{ridge}} = \arg \min_{\beta} & \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \\ \text{subject to} & \sum_{j=1}^p \beta_j^2 \leq t \end{aligned}$$

And there is a one-to-one correspondence between λ and t .

2.6.1 How to use ridge regression

Apparently, ridge regression is not equivariant under scaling of the inputs. For example, using OLS, measuring x in meters or in centimeters will not change the predictions, since the latter coefficients will just 100 times of the first. However, in ridge regression, the penalty term $\lambda \sum_{j=1}^p \beta_j^2$ will be affected by the scale of x_j . Which means, using centimeters instead of meters will make the penalty on β_j 10000 times larger, leading to different solutions. Thus, it's important to standardize the features (zero mean and unit variance) before applying ridge.

Usually, we calculate μ and σ from training set, and use them to standardize both training and test sets. Scaler can be regarded as part of the model, not data cleaning!

2.6.2 Equivalence of two forms

First, take a review of **KKT conditions**.

Consider a minimization problem with both equality and inequality constraints:

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{subject to } g_i(\mathbf{x}) &\leq 0, & i = 1, \dots, m \\ h_j(\mathbf{x}) &= 0, & j = 1, \dots, l \end{aligned}$$

And the lagrangian function is:

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^l \mu_j h_j(\mathbf{x})$$

If \mathbf{x}^* is a local minimum, then there exist multipliers $\lambda_i^* \geq 0$ and μ_j^* such that the following conditions hold:

- **Stationary**

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = 0$$

- **Primal feasibility**

$$\begin{aligned} g_i(\mathbf{x}^*) &\leq 0, & i = 1, \dots, m \\ h_j(\mathbf{x}^*) &= 0, & j = 1, \dots, l \end{aligned}$$

The gradient can be regarded as the force to push a particle, primal stationary means the force of $\partial f(\mathbf{x}^*)$ is balanced by a linear sum of forces from constraints.

- **Dual feasibility**

$$\lambda_i^* \geq 0, \quad i = 1, \dots, m$$

All the $\partial g_i(\mathbf{x}^*)$ forces must be one-sided, pointing inwards into the feasible set for \mathbf{x} .

- **Complementary slackness**

$$\lambda_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m$$

The force only activated when the particle is on the boundary of feasible set.

In ridge regression, we have:

$$\mathcal{L} = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \alpha(\|\beta\|_2^2 - t)$$

According to stationary condition:

$$\nabla_{\beta} \mathcal{L} = -2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta) + 2\alpha\beta = 0$$

On the other hand, solving the unconstrained form is

$$\nabla_{\beta} (\|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2) = -2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta) + 2\lambda\beta = 0$$

Thus, if we set $\lambda = \alpha$, the two forms are equivalent.

One step further, according to complementary slackness:

$$\alpha(\|\beta\|_2^2 - t) = 0$$

From unconstrained form, given λ , we can solve β , denote $\beta(\lambda)$, and define $t(\lambda) = \|\beta(\lambda)\|_2^2$. Then, $\beta(\lambda)$ is also the solution of constrained form with $t = t(\lambda)$ (apparently it's on the boundary, and the coefficient $\alpha = \lambda$).

From the constrained form, given t , we can also solve β , denote $\beta(t)$.

- If $\|\beta(t)\|_2^2 < t$, then according to complementary slackness, $\alpha = 0$, which means no penalty, $\lambda = 0$, back to OLS.
- If $\|\beta(t)\|_2^2 = t$, the boundary is effective, correspond to some $\alpha > 0$, and $\lambda = \alpha$.

2.6.3 Bayesian view

Assume prior $\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$, and likelihood $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I})$. Then the posterior is:

$$\begin{aligned} p(\beta|\mathbf{Y}) &\propto p(\mathbf{Y}|\beta)p(\beta) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{Y} - \mathbf{X}\beta\|_2^2\right) \exp\left(-\frac{1}{2\tau^2}\|\beta\|_2^2\right) \\ &\propto \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2}\|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \frac{1}{\tau^2}\|\beta\|_2^2\right)\right) \end{aligned}$$

Maximizing the posterior is equivalent to minimizing the negative log-posterior. The object is equivalent to ridge regression with $\lambda = \sigma^2/\tau^2$.

However, it seems that from bayesian view, we don't have a hard constrain on the size of β , just a prior that β is likely to be small. But the constrained form of ridge regression directly enforce $\|\beta\|_2^2 \leq t$. Is it contradictory?

No. Actually, the posterior distribution of β still have infinite support, meaning that β can still be large with small probability. Our MAP estimate of β is the **mode** of the posterior, it's just a **point estimate**. This point estimator can still satisfy the hard constraint $\|\beta\|_2^2 \leq t$ for some t .

In ridge regression, since the posterior is symmetric gaussian distribution, the mean and mode are same. However, in lasso and other priors, we usually use **mode**, **not mean** in common Bayes estimate.

2.6.4 SVD in ridge

The solution of ridge regression is:

$$\hat{\beta}^{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$$

When $\mathbf{X}^T \mathbf{X}$ is singular (not full rank), OLS estimator is not defined, and adding $\lambda \mathbf{I}$ ensures that $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is positive definite and invertible.

Take reduced SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, $r = \text{Rank}(\mathbf{X})$, where $\mathbf{U} \in \mathbb{R}^{N \times r}$, $\mathbf{D} \in \mathbb{R}^{r \times r}$, $\mathbf{V} \in \mathbb{R}^{p \times r}$. Then,

$$\begin{aligned} \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} &= \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T + \lambda \mathbf{I} \\ &= \mathbf{V} \mathbf{D}^T \mathbf{D} \mathbf{V}^T + \lambda \mathbf{V} \mathbf{V}^T \\ &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}) \mathbf{V}^T \end{aligned}$$

must be full rank and invertible.

It's also easy to see that in the case of orthogonal inputs (i.e., $\mathbf{X}^T \mathbf{X} = \mathbf{I}$), the ridge estimates are just a scaled version of the least squares estimates, that is

$$\hat{\beta}^{\text{ridge}} = \frac{1}{1 + \lambda} \hat{\beta}$$

Further more,

$$\begin{aligned} \mathbf{X} \beta^{\text{OLS}} &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{U} \mathbf{D} \mathbf{V}^T (\mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{Y} \\ &= \mathbf{U} \mathbf{D} \mathbf{D}^{-2} \mathbf{D}^T \mathbf{U}^T \mathbf{Y} \\ &= \mathbf{U} \mathbf{U}^T \mathbf{Y} \end{aligned}$$

and

$$\begin{aligned}
\mathbf{X}\beta^{\text{ridge}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{Y} \\
&= \mathbf{U}\mathbf{D}\mathbf{V}^T(\mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})\mathbf{V}^T)^{-1}\mathbf{V}\mathbf{D}^T\mathbf{U}^T\mathbf{Y} \\
&= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}^T\mathbf{U}^T\mathbf{Y} \\
&= \mathbf{U}\text{diag}\left(\frac{d_1^2}{d_1^2 + \lambda}, \frac{d_2^2}{d_2^2 + \lambda}, \dots, \frac{d_r^2}{d_r^2 + \lambda}\right)\mathbf{U}^T\mathbf{Y}
\end{aligned}$$

In both way, $\mathbf{U}^T\mathbf{Y}$ gets the coordinates of \mathbf{Y} in the basis of columns of \mathbf{U} (the principal components of \mathbf{X}), and then send back to original space using \mathbf{U} , but ridge regression shrinks the coordinates. It's easy to see the directions with larger variance have lower effect, and those with small variance shrink most.

Similar to OLS, the degree of freedom of ridge is defined as

$$\text{df}(\lambda) = \text{tr}(\mathbf{H}(\lambda)) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

But actually it's because the definition of df is

$$\text{df}(\hat{\mathbf{y}}) = \frac{1}{\sigma^2} \text{Cov}(\hat{y}_i, y_i)$$

indicating how much the prediction \hat{y}_i change if data y_i have a tiny change. If the change is big, then model is 'reciting' the data point, model is complex, taking lots of df. If the change is small, meaning model is smooth, little affected by single data point. Model is simple, df is low.

For those prediction is a linear transformation of \mathbf{y} (OLS, ridge), the df is the trace of hat matrix.

2.6.5 Revisit PCA

For a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, the mean is

$$\mu := \frac{1}{n} \mathbf{X}^T \mathbf{1} \in \mathbb{R}^p$$

The centered matrix is

$$\mathbf{X}_c := \mathbf{X} - \mathbf{1}\mu^T$$

centering is to put the center of data cloud in origin, otherwise PCs are not purely about variance, but also about the shift from origin.

The SVD of \mathbf{X}_c is $\mathbf{X}_c = \mathbf{U}\mathbf{D}\mathbf{V}^T$, and since it's centered, the covariance is

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}_c^T \mathbf{X}_c = \frac{1}{n-1} \mathbf{V}\mathbf{D}^2\mathbf{V}^T$$

apparently, each column of \mathbf{V} is the eigenvector of \mathbf{S} : $v_k = \mathbf{V}_{:,k} \in \mathbb{R}^p$, and project origin p dim feature into one dim.

Thus, to get a k dim representation, we need to take first k column of \mathbf{V} , which is an encoder:

$$\mathbf{Z}_k = \mathbf{X}_c \mathbf{V}_k = \mathbf{U}_k \mathbf{\Sigma}_k$$

The first principal component direction $z_1 = \mathbf{X}_c v_1$ has **largest** sample variance among all normalized linear combinations of the columns of \mathbf{X}_c , the variance is

$$\begin{aligned}\text{Var}(z_1) &= \frac{1}{n-1} z_1^T z_1 = \frac{1}{n-1} (\mathbf{X}_c v_1)^T (\mathbf{X}_c v_1) \\ &= \frac{1}{n-1} v_1^T \mathbf{X}_c^T \mathbf{X}_c v_1 \\ &= v_1^T \mathbf{S} v_1 \\ &= v_1^T \lambda_1 v_1 \\ &= \lambda_1 = \frac{\sigma_1^2}{n-1}\end{aligned}$$

So the variance of the first PC scores is exactly the top eigenvalue of the covariance matrix.

And if we need to reconstruct the origin matrix, just decode in the same way:

$$\hat{\mathbf{X}}_c = \mathbf{Z}_k \mathbf{V}_k^T, \quad \hat{\mathbf{X}} = \hat{\mathbf{X}}_c + \mathbf{1} \mu^T$$

for a single sample:

$$\hat{x} = \mathbf{V}_k z = \mathbf{V}_k (\mathbf{V}_k^T z)$$

Thus, \mathbf{V} is just orthogonal basis of the feature space. It maps centered data in to its space, and get the coordinates.

From an auto encoder's view, it is

$$\min_{\mathbf{W} \in \mathbb{R}^{p \times k}} \|\mathbf{X}_c - \mathbf{X}_c \mathbf{W} \mathbf{W}^T\|_F^2 \quad \text{s.t.} \quad \mathbf{W}^T \mathbf{W} = \mathbf{I}$$

It seems we are talking about \mathbf{V} all the time. What about \mathbf{U} ? Well, he is the basis of the sample space. Since we don't care about sample space that much, we just don't talk about it that much.

2.7 LASSO

2.7.1 Subgradient and subdifferential

For a convex function $f : \mathbb{R} \rightarrow \mathbb{R}$, a number g is a subgradient of f at x_0 if for all x ,

$$f(x) \geq f(x_0) + g \cdot (x - x_0)$$

So the line with slope g touching $(x_0, f(x_0))$ lies below the whole function.

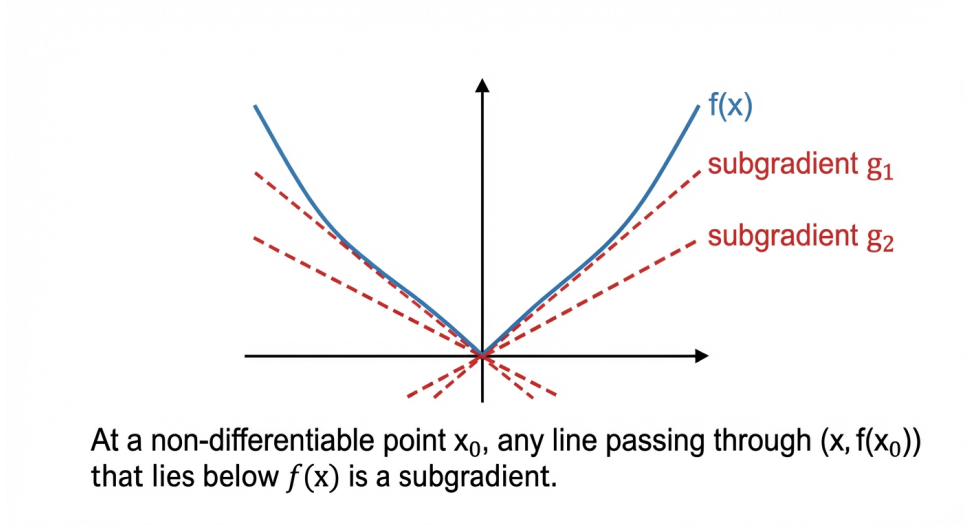


Figure 1: Illustration of subgradient.

The set of all such slopes is the subdifferential:

$$\partial f(x_0) = \{g : f(x) \geq f(x_0) + g \cdot (x - x_0), \forall x\}$$

Thus, for $f(x) = |x|$,

$$\partial|x| = \begin{cases} \{1\} & \text{if } x > 0 \\ \{-1\} & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \end{cases}$$

When $x = 0$, the subdifferential is an interval!

Common way to compute Lasso is ISAT in 4.4.1.

2.7.2 Why shrink to zero

Consider LASSO objective function

$$\min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

the solution $\hat{\beta}$ satisfies

$$\mathbf{0} \in \partial J(\hat{\beta}) \in \mathbb{R}^p$$

take one feature, and define $\mathbf{r} = \mathbf{Y} - \mathbf{X}\hat{\beta}$

$$0 \in -X_j^T \mathbf{r} + \lambda \partial \hat{\beta}_j$$

which means, if $X_j^T \mathbf{r} \in [-\lambda, \lambda]$, $\hat{\beta}_j$ will be shrink to zero.

What does a low value of $\mathbf{x}_j^T \mathbf{r}$ actually indicate? It signifies that the feature has a low correlation with the current residual. In other words, this feature is not useful for reducing the prediction error, even if it contains 'new' or unique information.

We must distinguish this from **feature-to-feature correlation** (multicollinearity):

- **Redundancy (High Feature Correlation):** Consider predicting weight using height. Height in 'cm' and 'inches' are perfectly correlated. Lasso chooses only one because adding the second provides no *additional* help in reducing the residual.
- **Irrelevance (Low Residual Correlation):** Consider the feature 'eating carrots'. This is uncorrelated with height (it brings independent information), but it has no relationship with weight. Since it cannot lower the residual ($\mathbf{x}_{carrots}^T \mathbf{r} \approx 0$), Lasso abandons it.

2.7.3 Soft thresholding

When design matrix \mathbf{X} is orthogonal, i.e. $\mathbf{X}^T \mathbf{X} = \mathbf{I}$, the OLS solution become $\hat{\boldsymbol{\beta}}^{\text{OLS}} = \mathbf{X}^T \mathbf{Y}$, and LASSO becomes

$$\begin{aligned}\hat{\boldsymbol{\beta}}^{\text{Lasso}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\mathbf{X}^T \mathbf{Y} - \boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\hat{\boldsymbol{\beta}}^{\text{OLS}} - \boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1\end{aligned}$$

Let the gradient to zero,

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\beta}} &= \boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{\text{OLS}} + \lambda \cdot \operatorname{sign}(\boldsymbol{\beta}) = 0 \\ \Rightarrow \quad \hat{\boldsymbol{\beta}}^{\text{Lasso}} &= \hat{\boldsymbol{\beta}}^{\text{OLS}} - \lambda \cdot \operatorname{sign}(\hat{\boldsymbol{\beta}}^{\text{Lasso}}) \\ &= \operatorname{sign}(\hat{\boldsymbol{\beta}}^{\text{OLS}}) \max\{0, |\hat{\boldsymbol{\beta}}^{\text{OLS}}| - \lambda\} \\ &= S_\lambda(\hat{\boldsymbol{\beta}}^{\text{OLS}})\end{aligned}$$

Here, S_λ is the soft-thresholding operator.

3 Chap 4: Linear methods for classification

Generally, there are two ways to model classification problem: hard boundary $\delta_k(x)$ or posterior probabilities $\Pr(G = k \mid X = x)$. If either $\delta_k(x)$ or $\Pr(G = k \mid X = x)$ are linear in x , then the decision boundary is linear.

3.1 Linear Regression of an Indicator Matrix

Let \mathbf{Y} be the $N \times K$ indicator matrix where each row has exactly one 1. Let \mathbf{X} be the $N \times (p+1)$ design matrix, including a leading column of ones for the intercept.

We fit a linear model to **each column** of \mathbf{Y} simultaneously using Least Squares. The objective is to minimize the Frobenius norm of the residual matrix:

$$\min_{\mathbf{B}} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2$$

The solution $\hat{\mathbf{B}}$ is given by the normal equation:

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

For a new input vector $\mathbf{x} = (1, x_1, \dots, x_p)^T$, the vector of fitted values (class scores) is:

$$\hat{\mathbf{f}}(\mathbf{x})^T = \mathbf{x}^T \hat{\mathbf{B}} = \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Proof: Why $\sum \hat{f}_k(\mathbf{x}) = 1$

We want to calculate the sum of the predicted values across all K classes. Let $\mathbf{1}_K$ be a column vector of K ones, and $\mathbf{1}_N$ be a column vector of N ones.

The sum can be written as the dot product of the prediction vector and $\mathbf{1}_K$:

$$\sum_{k=1}^K \hat{f}_k(\mathbf{x}) = \hat{\mathbf{f}}(\mathbf{x})^T \mathbf{1}_K = \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underbrace{(\mathbf{Y} \mathbf{1}_K)}_{\text{step (a)}}$$

Step (a): Since each row of the indicator matrix \mathbf{Y} sums to exactly 1 (one-hot encoding), multiplying \mathbf{Y} by $\mathbf{1}_K$ results in a vector of ones:

$$\mathbf{Y} \mathbf{1}_K = \mathbf{1}_N$$

Substituting this back:

$$\sum_{k=1}^K \hat{f}_k(\mathbf{x}) = \mathbf{x}^T \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{1}_N}_{\text{step (b)}}$$

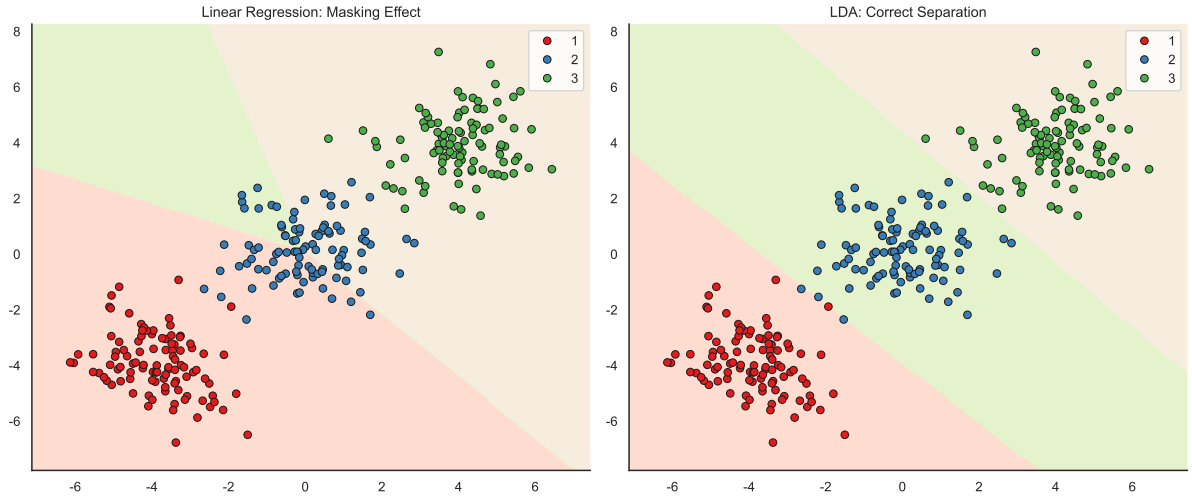
Step (b): The term $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{1}_N$ represents the OLS regression coefficients obtained when regressing the vector $\mathbf{1}_N$ onto \mathbf{X} . Since \mathbf{X} contains a column of ones (the intercept), the constant vector $\mathbf{1}_N$ lies perfectly in the column space of \mathbf{X} . The coefficients that strictly reconstruct $\mathbf{1}_N$ are 1 for the intercept and 0 for all other features. Let $\mathbf{e}_1 = (1, 0, \dots, 0)^T$. Then:

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{1}_N = \mathbf{e}_1$$

Conclusion: Since the input vector is $\mathbf{x} = (1, x_1, \dots, x_p)^T$, its dot product with \mathbf{e}_1 picks out the first element (the intercept):

$$\sum_{k=1}^K \hat{f}_k(\mathbf{x}) = \mathbf{x}^T \mathbf{e}_1 = 1$$

Masking effect



4 Optimization methods

I read [1] for convex optimization methods.

4.1 Norms of matrices

Mathematically, a function $f : \mathbb{K}^{m \times n} \rightarrow \mathbb{R}$ is a matrix norm if for all matrices $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{m \times n}$ and scalars α , it satisfies:

- Non-negativity: $\|\mathbf{A}\| \geq 0$ (and $\|\mathbf{A}\| = 0 \iff \mathbf{A} = \mathbf{0}$).
- Scalar Property: $\|\alpha\mathbf{A}\| = |\alpha|\|\mathbf{A}\|$.
- Triangle Inequality: $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$.

Most useful matrix norms also satisfy the submultiplicative property:

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\|\|\mathbf{B}\|$$

Two kinds of norms: Induced (operator) norms and Frobenius norm (not induced).

4.1.1 Induced norm

An **induced norm** is defined by how much the matrix \mathbf{A} stretches a vector \mathbf{x} under a specific vector p -norm.

$$\|\mathbf{A}\|_p = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{Ax}\|_p$$

The "Big Three" Induced Norms

- L_1 Norm (Maximum Column Sum):

$$\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$$

- L_∞ Norm (Maximum Row Sum):

$$\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$$

- L_2 Norm (Spectral Norm):

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^H \mathbf{A})} = \sigma_{\max}(\mathbf{A})$$

where σ_{\max} is the largest singular value of \mathbf{A} .

Here are proofs

- **L_∞ Norm (Maximum Row Sum)**

Step 1: Upper Bound. For any vector \mathbf{x} such that $\|\mathbf{x}\|_\infty = 1$, we have $|x_j| \leq 1$ for all j . For any row i :

$$|(\mathbf{Ax})_i| = \left| \sum_{j=1}^n a_{ij}x_j \right| \leq \sum_{j=1}^n |a_{ij}||x_j| \leq \sum_{j=1}^n |a_{ij}|$$

Thus, $\|\mathbf{Ax}\|_\infty = \max_i |(\mathbf{Ax})_i| \leq \max_i \sum_{j=1}^n |a_{ij}|$.

Step 2: Attainability. Let k be the row index where $\sum_j |a_{kj}|$ is maximized. Define \mathbf{x} such that $x_j = \text{sgn}(a_{kj})$ if $a_{kj} \neq 0$ and $x_j = 1$ otherwise. Then $\|\mathbf{x}\|_\infty = 1$ and $|(\mathbf{Ax})_k| = \sum_j |a_{kj}|$, achieving the bound.

- **L_1 Norm (Maximum Column Sum)**

Step 1: Upper Bound. For $\|\mathbf{x}\|_1 = 1$, we expand the L_1 norm of the product:

$$\|\mathbf{Ax}\|_1 = \sum_{i=1}^m \left| \sum_{j=1}^n a_{ij}x_j \right| \leq \sum_{i=1}^m \sum_{j=1}^n |a_{ij}||x_j| = \sum_{j=1}^n |x_j| \left(\sum_{i=1}^m |a_{ij}| \right)$$

here, for finite and non-negative summation, we can always switch the order.

Since $\sum_i |a_{ij}| \leq \max_j \sum_i |a_{ij}|$, we have:

$$\|\mathbf{Ax}\|_1 \leq \left(\max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \right) \sum_{j=1}^n |x_j| = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$$

Step 2: Attainability. Let k be the index of the column with the maximum absolute sum. Set $\mathbf{x} = \mathbf{e}_k$ (the k -th canonical basis vector). Then $\|\mathbf{x}\|_1 = 1$ and $\|\mathbf{Ax}\|_1 = \sum_i |a_{ik}|$, achieving the bound.

- **L_2 Norm (Spectral Norm)**

Step 1: Upper Bound. Consider the square of the norm $\|\mathbf{Ax}\|_2^2 = \mathbf{x}^H \mathbf{A}^H \mathbf{Ax}$. Let $\mathbf{M} = \mathbf{A}^H \mathbf{A}$, which is Hermitian and positive semi-definite. By the spectral theorem, for any \mathbf{x} with $\|\mathbf{x}\|_2 = 1$:

$$\mathbf{x}^H \mathbf{M} \mathbf{x} \leq \lambda_{\max}(\mathbf{M}) \|\mathbf{x}\|_2^2 = \lambda_{\max}(\mathbf{A}^H \mathbf{A})$$

Taking the square root, we find $\|\mathbf{Ax}\|_2 \leq \sqrt{\lambda_{\max}(\mathbf{A}^H \mathbf{A})}$.

Step 2: Attainability. Choose \mathbf{x} to be the normalized eigenvector of $\mathbf{A}^H \mathbf{A}$ corresponding to λ_{\max} . Then $\mathbf{x}^H \mathbf{A}^H \mathbf{Ax} = \mathbf{x}^H (\lambda_{\max} \mathbf{x}) = \lambda_{\max} \|\mathbf{x}\|_2^2 = \lambda_{\max}$, achieving the bound.

4.1.2 Frobenius norm

This is the most common entry-wise norm, analogous to the Euclidean distance for vectors.

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{tr}(\mathbf{A}^H \mathbf{A})}$$

Where tr is the trace and \mathbf{A}^H is the conjugate transpose. It is also equal to the square root of the sum of the squares of the singular values:

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^{\min(m,n)} \sigma_i^2}$$

4.1.3 Schatten norm

Schatten norms apply a p -norm to the singular values σ_i of the matrix. Let $\sigma(\mathbf{A})$ be the vector of singular values:

$$\|\mathbf{A}\|_p = \left(\sum_{i=1}^{\min(m,n)} \sigma_i^p \right)^{1/p}$$

- $p = 1$ (Nuclear Norm): $\|\mathbf{A}\|_* = \sum \sigma_i$. Often used in machine learning for rank minimization.
- $p = 2$ (Frobenius Norm): $\|\mathbf{A}\|_F = \sqrt{\sum \sigma_i^2}$.
- $p = \infty$ (Spectral Norm): $\|\mathbf{A}\|_2 = \sigma_{\max}$.

4.2 Gradient descent

Assume f is convex and differentiable, with $\text{dom}(f) = \mathbb{R}^n$, and additionally that ∇f is Lipschitz continuous with constant $L > 0$,

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|_2, \quad \forall x, y$$

when the stepsize $t \leq 1/L$, we have

$$f(x^{(k)}) - f^* \leq \frac{\|x^{(0)} - x^*\|^2}{2tk}$$

indicating a convergence rate $\mathcal{O}(1/k)$. In another word, to achieve precision ϵ , needs about $\mathcal{O}(1/\epsilon)$ times of iteration.

4.2.1 Descent Lemma

Express the difference between two points

$$f(\mathbf{y}) = f(\mathbf{x}) + \int_0^1 \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x}))^T (\mathbf{y} - \mathbf{x}) d\tau$$

then introduce the gradient at \mathbf{x}

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \int_0^1 (\nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}))^T (\mathbf{y} - \mathbf{x}) d\tau$$

apply Cauchy-Schwarz Inequality,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \int_0^1 \|\nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\| d\tau$$

Using the L -Lipschitz property, we know $\|\nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \leq L\|\tau(\mathbf{y} - \mathbf{x})\| = L\tau\|\mathbf{y} - \mathbf{x}\|$. Substituting this back:

$$\begin{aligned} f(\mathbf{y}) &\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \int_0^1 L\tau\|\mathbf{y} - \mathbf{x}\|^2 d\tau \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2}L\|\mathbf{y} - \mathbf{x}\|^2 \end{aligned}$$

4.2.2 Convergence rate

In gradient descent update $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t\nabla f(\mathbf{x}^{(k)})$ with step size $t \leq 1/L$:

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &\leq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + \frac{L}{2} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^2 \\ &= f(\mathbf{x}^{(k)}) - t\|\nabla f(\mathbf{x}^{(k)})\|^2 + \frac{Lt^2}{2} \|\nabla f(\mathbf{x}^{(k)})\|^2 \\ &= f(\mathbf{x}^{(k)}) + \left(\frac{Lt^2}{2} - t\right) \|\nabla f(\mathbf{x}^{(k)})\|^2 \end{aligned}$$

Since $t \leq 1/L$, we have $-t + \frac{Lt^2}{2} \leq -\frac{t}{2}$. Thus:

$$f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \frac{t}{2} \|\nabla f(\mathbf{x}^{(k)})\|^2$$

which ensures **descent**.

By the convexity of f , we have $f(\mathbf{x}^*) \geq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^* - \mathbf{x}^{(k)})$, which rearranges to:

$$f(\mathbf{x}^{(k)}) - f^* \leq \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^{(k)} - \mathbf{x}^*)$$

plus the previous one, and get

$$f(\mathbf{x}^{(k+1)}) - f^* \leq \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^{(k)} - \mathbf{x}^*) - \frac{t}{2} \|\nabla f(\mathbf{x}^{(k)})\|^2$$

What's next? We want to know the relationship between the descent of $f(\mathbf{x}^{(k)} - f^*)$, and the descent of $\mathbf{x}^{(k)} - \mathbf{x}^*$, so consider

$$\begin{aligned} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 &= \|(\mathbf{x}^{(k)} - \mathbf{x}^*) - t\nabla f(\mathbf{x}^{(k)})\|^2 \\ &= \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - 2t\nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^{(k)} - \mathbf{x}^*) + t^2 \|\nabla f(\mathbf{x}^{(k)})\|^2 \end{aligned}$$

thus,

$$f(\mathbf{x}^{(k+1)}) - f^* \leq \frac{1}{2t} \left(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 \right)$$

Summing from $i = 0$ to $k - 1$:

$$\sum_{i=0}^{k-1} (f(\mathbf{x}^{(i+1)}) - f^*) \leq \frac{1}{2t} \left(\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 \right) \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{2t}$$

and we know $k(f(\mathbf{x}^{(k)}) - f^*) \leq \sum_{i=0}^{k-1} (f(\mathbf{x}^{(i+1)}) - f^*)$, we get

$$f(\mathbf{x}^{(k)}) - f^* \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{2tk}$$

4.2.3 Another view of GD

If we approximate the function around \mathbf{x} by replacing the Hessian $\nabla^2 f(\mathbf{x})$ with $\frac{1}{t}\mathbf{I}$:

$$\mathbf{x}^+ = \arg \min_{\mathbf{z}} \left(f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{z} - \mathbf{x}) + \frac{1}{2t} \|\mathbf{z} - \mathbf{x}\|_2^2 \right)$$

where t is the step size. Surprisingly, taking the derivative of the right-hand side with respect to \mathbf{z} and setting it to zero gives the gradient descent update:

$$\nabla f(\mathbf{x}) + \frac{1}{t}(\mathbf{z} - \mathbf{x}) = 0 \implies \mathbf{x}^+ = \mathbf{x} - t\nabla f(\mathbf{x})$$

From this view, the standard gradient descent assume the Hessian matrix is same.

4.3 Subgradient descent

Here we assume function is G Lipschitz continuous: $|f(\mathbf{x}) - f(\mathbf{y})| \leq G\|\mathbf{x} - \mathbf{y}\|_2$, which bounds the norm of subgradient $\|\mathbf{g}\|_2 \leq G$.

The subgradient update rule is $\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - t_k \mathbf{g}^{(k-1)}$, where $\mathbf{g}^{(k-1)} \in \partial f(\mathbf{x}^{(k-1)})$.

Analyze the squared distance to the optimal point \mathbf{x}^* :

$$\begin{aligned}\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 &= \|\mathbf{x}^{(k-1)} - t_k \mathbf{g}^{(k-1)} - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|_2^2 - 2t_k (\mathbf{g}^{(k-1)})^T (\mathbf{x}^{(k-1)} - \mathbf{x}^*) + t_k^2 \|\mathbf{g}^{(k-1)}\|_2^2\end{aligned}$$

From the definition of a subgradient: $(\mathbf{g}^{(k-1)})^T (\mathbf{x}^{(k-1)} - \mathbf{x}^*) \geq f(\mathbf{x}^{(k-1)}) - f(\mathbf{x}^*)$. Substituting this and $\|\mathbf{g}\|_2 \leq G$:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}^{(k-1)} - \mathbf{x}^*\|_2^2 - 2t_k (f(\mathbf{x}^{(k-1)}) - f^*) + t_k^2 G^2 \quad (*)$$

Summing from $i = 1$ to k and noting $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 \geq 0$:

$$0 \leq R^2 - 2 \sum_{i=1}^k t_i (f(\mathbf{x}^{(i-1)}) - f^*) + G^2 \sum_{i=1}^k t_i^2$$

where $R = \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2$. Rearranging for the best iterate $f(\mathbf{x}_{best}^{(k)}) = \min_{i=0, \dots, k} f(\mathbf{x}^{(i)})$:

$$f(\mathbf{x}_{best}^{(k)}) - f^* \leq \frac{R^2 + G^2 \sum_{i=1}^k t_i^2}{2 \sum_{i=1}^k t_i}$$

Thus, if $\sum_{i=1}^k t_i^2 < \infty$ and $\sum_{i=1}^k t_i = \infty$, subgradient can converge to optimal.

For a fixed step size $t_i = t$:

$$f(\mathbf{x}_{best}^{(k)}) - f^* \leq \frac{R^2}{2kt} + \frac{G^2 t}{2}$$

To achieve an error of ϵ , we set $t = \epsilon/G^2$. This requires:

$$k \geq \frac{R^2 G^2}{\epsilon^2}$$

Thus, the subgradient method has a convergence rate of $\mathcal{O}(1/\epsilon^2)$, which is slower than the $\mathcal{O}(1/\epsilon)$ rate of gradient descent.

When optimal value f^* is known, minimizing equation (*) get Polyak step sizes:

$$t_k = \frac{f(\mathbf{x}^{(k-1)}) - f^*}{\|\mathbf{g}^{(k-1)}\|_2^2}$$

4.4 Proximal Gradient Descent

Although subgradient methods can solve non-differentiable optimization problems, it's slow $\mathcal{O}(1/\epsilon^2)$. Thus, we try write it as a composite function

$$f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$$

where g is convex and differentiable, h is convex but not necessarily differentiable.

For smooth g , use quadratic approximation, and for non-smooth h , just keep it.

$$\begin{aligned}\mathbf{x}^+ &= \underset{\mathbf{z}}{\operatorname{argmin}} \left(g(\mathbf{x}) + \nabla g(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) + \frac{1}{2t} \|\mathbf{z} - \mathbf{x}\|_2^2 \right) + h(\mathbf{z}) \\ &= \underset{\mathbf{z}}{\operatorname{argmin}} \left(\frac{1}{2t} \|\mathbf{z} - (\mathbf{x} - t\nabla g(\mathbf{x}))\|_2^2 + h(\mathbf{z}) \right)\end{aligned}$$

The first term make \mathbf{z} stay close to the gradient update of g , while also make h small (second term).

Define **proximal mapping**:

$$\operatorname{prox}_{h,t}(\mathbf{v}) = \underset{\mathbf{z}}{\operatorname{argmin}} \left(\frac{1}{2t} \|\mathbf{z} - \mathbf{v}\|_2^2 + h(\mathbf{z}) \right)$$

then the procedure of proximal gradient descent is, for $k = 1, 2, 3, \dots$,

$$\mathbf{v} = \mathbf{x}^{(k-1)} - t_k \nabla g(\mathbf{x}^{(k-1)})$$

$$\mathbf{x}^{(k)} = \operatorname{prox}_{h,t_k}(\mathbf{v})$$

The name of *proximal* comes from $\|\mathbf{z} - \mathbf{v}\|$ term, it makes the new solution \mathbf{z} near input point \mathbf{v}

Proximal gradient descent have same convergence rate as gradient descent, but need to consider prox cost...

4.4.1 ISTA

Iterative Shrinkage-Thresholding Algorithm (ISTA) is used for Lasso [2.7](#).

$$\min \frac{1}{2} \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

$$\operatorname{prox}_{\lambda\|\cdot\|_1,t}(\mathbf{v}) = \underset{\mathbf{z}}{\operatorname{argmin}} \left(\frac{1}{2t} \|\mathbf{z} - \mathbf{v}\|_2^2 + \lambda \|\mathbf{z}\|_1 \right)$$

The solution is soft-thresholding operator, denote $S_{\lambda t}(\mathbf{v})$

$$[S_{\lambda t}(\mathbf{v})]_i = \begin{cases} v_i - \lambda t & \text{if } v_i > \lambda t \\ 0 & \text{if } |v_i| \leq \lambda t \\ v_i + \lambda t & \text{if } v_i < -\lambda t \end{cases}$$

In practice, iteratively do

- Residual gradient update: $\mathbf{v} = \boldsymbol{\beta} + tX^T(\mathbf{y} - X\boldsymbol{\beta})$
- Soft-thresholding truncate: $\boldsymbol{\beta}^+ = S_{\lambda t}(\mathbf{v})$

4.5 Stochastic gradient descent

References

- [1] Ryan Tibshirani. *Convex Optimization (10-725/36-725)*. Carnegie Mellon University course website. Accessed: 2026-01-03. 2019. URL: <https://www.stat.cmu.edu/~ryantibs/convexopt/>.
- [2] Wikipedia contributors. *Gauss–Markov theorem* — *Wikipedia, The Free Encyclopedia*. Accessed: 2025-12-29. 2025. URL: https://en.wikipedia.org/wiki/Gauss%E2%80%93Markov_theorem.