# Benign Overfitting in Linear Regression

Peter L. Bartlett[a,b,1], Philip M. Long[c] iD, Gábor Lugosi[d,e,f] iD, and Alexander Tsigler[a]

statistical learning theory | overfitting | linear regression | interpolation

# Contents

1. Research Background

2. Previous Studies

3. Definitions and Notation

4. Theorems and Conclusions

   Th1, Th2

5. Proof（Th1）

   Le1~Le11

6. Research Prospects and Future

# Benign overfitting in linear regression

Peter L. Bartlett[a,b,1], Philip M. Long[c], Gábor Lugosi[d,e,f], and Alexander Tsigler[a]

[a]Department of Statistics, University of California, Berkeley, CA 94720-3860; [b]Computer Science Division, University of California, Berkeley, CA 94720-1776; [c]Google Brain, Mountain View, CA 94043; [d]Economics and Business, Pompeu Fabra University, 08005 Barcelona, Spain; [e]Institució Catalana de Recerca i Estudis Avançats, Passeig, Lluís Companys 23, 08010 Barcelona, Spain; and [f]Barcelona Graduate School of Economics, 08005 Barcelona, Spain

The phenomenon of benign overfitting is one of the key mysteries uncovered by deep learning methodology: deep neural networks seem to predict well, even with a perfect fit to noisy training data. Motivated by this phenomenon, we consider when a perfect fit to training data in linear regression is compatible with accurate prediction. We give a characterization of linear regression problems for which the minimum norm interpolating prediction rule has near-optimal prediction accuracy. The characterization is in terms of two notions of the effective rank of the data covariance. It shows that overparameterization is essential for benign overfitting in this setting: the number of directions in parameter space that are unimportant for prediction must significantly exceed the sample size. By studying examples of data covariance properties that this characterization shows are required for benign overfitting, we find an important role for finite-dimensional data: the accuracy of the minimum norm interpolating prediction rule approaches the best possible accuracy for a much narrower range of properties of the data distribution when the data lie in an infinite-dimensional space vs. when the data lie in a finite-dimensional space with dimension that grows faster than the sample size.

statistical learning theory | overfitting | linear regression | interpolation

Deep learning methodology has revealed a surprising statistical phenomenon: overfitting can perform well. The classical perspective in statistical learning theory is that there should be a tradeoff between the fit to the training data and the complexity of the prediction rule. Whether complexity is measured in terms of the number of parameters, the number of nonzero parameters in a high-dimensional setting, the number of neighbors averaged in a nearest neighbor estimator, the scale of an estimate in a reproducing kernel Hilbert space, or the bandwidth of a kernel smoother, this tradeoff has been ubiquitous in statistical learning theory. Deep learning seems to operate outside the regime where results of this kind are informative since deep neural networks can perform well even with a perfect fit to the training data.

As one example of this phenomenon, consider the experiment illustrated in figure 1C in ref. 1: standard deep network architectures and stochastic gradient algorithms, run until they perfectly fit a standard image classification training set, give respectable prediction performance, even when significant levels of label noise are introduced. The deep networks in the experiments reported in ref. 1 achieved essentially zero cross-entropy loss on the training data. In statistics and machine learning textbooks, an estimate that fits every training example perfectly is often presented as an illustration of overfitting ["…interpolating fits…[are] unlikely to predict future data well at all" (ref. 2, p. 37)]. Thus, to arrive at a scientific understanding of the success of deep learning methods, it is a central challenge to understand the performance of prediction rules that fit the training data perfectly.

In this paper, we consider perhaps the simplest setting where we might hope to witness this phenomenon: linear regression. That is, we consider quadratic loss and linear prediction rules, and we assume that the dimension of the parameter space is large

enough that a perfect fit is guaranteed. We consider data in an infinite-dimensional space (a separable Hilbert space), but our results apply to a finite-dimensional subspace as a special case. There is an ideal value of the parameters, $\theta^*$, corresponding to the linear prediction rule that minimizes the expected quadratic loss. We ask when it is possible to fit the data exactly and still compete with the prediction accuracy of $\theta^*$. Since we require more parameters than the sample size in order to fit exactly, the solution might be underdetermined, and therefore, there might be many interpolating solutions. We consider the most natural: choose the parameter vector $\hat{\theta}$ with the smallest norm among all vectors that gives perfect predictions on the training sample. (This corresponds to using the pseudoinverse to solve the normal equations; see below.) We ask when it is possible to overfit in this way—and embed all of the noise of the labels into the parameter estimate $\hat{\theta}$—without harming prediction accuracy.

Our main result is a finite sample characterization of when overfitting is benign in this setting. The linear regression problem depends on the optimal parameters $\theta^*$ and the covariance $\Sigma$ of the covariates $x$. The properties of $\Sigma$ turn out to be crucial since the magnitude of the variance in different directions determines both how the label noise gets distributed across the parameter space and how errors in parameter estimation in different directions in parameter space affect prediction accuracy. There is a classical decomposition of the excess prediction error into two terms. The first is rather standard: provided that the scale of the problem (that is, the sum of the eigenvalues of $\Sigma$) is small compared with the sample size $n$, the contribution to $\hat{\theta}$ that we can view as coming from $\theta^*$ is not too distorted. The second term is more interesting since it reflects the impact of the noise in the labels on prediction accuracy. We show that this part is small if and only if the effective rank of $\Sigma$ in the subspace corresponding to low-variance directions is large compared with $n$. This necessary and sufficient condition of a large effective rank can be viewed as a property of significant overparameterization: fitting the training data exactly but with near-optimal prediction accuracy occurs if and only if there are many low-variance (and

STATISTICS

# 1. Research Background

- A surprising phenomenon in seep learning methodology: **overfitting can perform well**.

- Whether complexity is measured in terms of :

- the number of **parameters**,

- the number of **nonzero parameters** in a high-dimensional setting,

- the number of **neighbors** averaged in a **nearest neighbor estimator**,

- the scale of an estimate in a **reproducing kernel Hilbert space**,

- the bandwidth of a **kernel smoother**,

- this tradeoff has been ubiquitous in **statistical learning theory**.

# 1. Research Background

- **Deep learning** seems to operate outside the regime where results of this kind are informative since deep neural networks can **perform well** even with **a perfect fit to the training data**.

- To arrive at a scientific understanding of the success of deep learning methods, it is a central challenge to **understand the performance of prediction rules** that fit the training data perfectly.

# 1. Research Background

- The simplest setting: **Linear Regression**.


- Consider **quadratic loss** and **linear prediction rules**

- Assume that the **dimension** of the parameter space is **large** enough

- Consider data in an **infinite-dimensional** space


- We ask when it is possible to **fit the data exactly** and still compete with the prediction accuracy of $\theta^*$

# 1. Research Background

- In an **infinite-dimensional setting**, benign overfitting occurs only for **a narrow range of decay rates** of the eigenvalues.

- On the other hand, it occurs with **any suitably slowly decaying eigenvalue** sequence **in a finite-dimensional space** with **dimension** that grows **faster** than the sample size.

# 2.Previous Studies

**1）. Initial Observations & Experimental Studies**

Interpolating prediction rules (fitting noisy data exactly) emerged as a key mystery since 2017 Simons Institute program.

Belkin et al. (3) experimentally showed RKHS interpolants achieve high accuracy despite violating classical generalization bounds.

**2）. Interpolating Decision Rules & Kernel Methods**

Belkin et al. (4) proposed "simplicial interpolation" with asymptotic consistency in high dimensions.

Belkin et al. (5) studied singular-kernel smoothing methods achieving optimal nonparametric rates while interpolating (building on (6)).

Liang & Rakhlin (7) proved minimum-norm kernel regression (with nonlinear inner-product kernels) has good accuracy under specific sample conditions.

**3）. Parameter Space Dimension & Excess Risk**

Belkin et al. (8) experimentally analyzed excess risk as a function of parameter space dimension in linear/nonlinear models.

# 2.Previous Studies

**4）. Subsequent Linear Model Analyses**

Belkin et al. (11) derived excess risk for linear models (sparse parameters/Fourier features).

Hastie et al. (12) studied linear regression asymptotically ($n, p \to \infty, p/n \to \gamma$)):

Assumed convergence of spectral distribution of $\Sigma$, used random matrix theory to characterize excess prediction error;

Examined effects of noise variance, eigenvalue distribution, and $\gamma$;

Extended to models with random nonlinear features.

**5）. Contrast with Present Work**

Key contributions of this work:

Provides finite-sample upper/lower bounds for excess prediction error under arbitrary covariances and dimensions;

Requires no asymptotic assumptions or specific data distributions

# 3. Definitions and Notations

***Definition 1 (Linear Regression):*** A linear regression problem in a separable Hilbert space $\mathbb{H}$ is defined by a random covariate vector $x \in \mathbb{H}$ and outcome $y \in \mathbb{R}$. We define

1) the covariance operator $\Sigma = \mathbb{E}[xx^\top]$ and
2) the optimal parameter vector $\theta^* \in \mathbb{H}$, satisfying $\mathbb{E}(y - x^\top \theta^*)^2 = \min_\theta \mathbb{E}(y - x^\top \theta)^2$.

We assume that

1) $x$ and $y$ are mean zero;
2) $x = V \Lambda^{1/2} z$, where $\Sigma = V \Lambda V^\top$ is the spectral decomposition of $\Sigma$ and $z$ has components that are independent $\sigma_x^2$ sub-Gaussian with $\sigma_x$ a positive constant: that is, for all $\lambda \in \mathbb{H}$,

$$\mathbb{E}[\exp(\lambda^\top z)] \leq \exp(\sigma_x^2 \|\lambda\|^2 / 2),$$

where $\| \cdot \|$ is the norm in the Hilbert space $\mathbb{H}$;

3) the conditional noise variance is bounded below by some constant $\sigma^2$,

$$\mathbb{E}\left[(y - x^\top \theta^*)^2 \big| x \right] \geq \sigma^2;$$

4) $y - x^\top \theta^*$ is $\sigma_y^2$ sub-Gaussian conditionally on $x$: that is, for all $\lambda \in \mathbb{R}$,

$$\mathbb{E}[\exp(\lambda(y - x^\top \theta^*))|x] \leq \exp(\sigma_y^2 \lambda^2 / 2)$$

(note that this implies $\mathbb{E}[y|x] = x^\top \theta^*$); and
5) almost surely, the projection of the data $X$ on the space orthogonal to any eigenvector of $\Sigma$ spans a space of dimension $n$.

Given a training sample $(x_1, y_1), \ldots, (x_n, y_n)$ of $n$ independent pairs with the same distribution as $(x, y)$, an estimator returns a parameter estimate $\theta \in \mathbb{H}$. The excess risk of the estimator is defined as

$$R(\theta) := \mathbb{E}_{x,y}\left[\left(y - x^\top \theta\right)^2 - \left(y - x^\top \theta^*\right)^2\right],$$

# 3. Definitions and Notations

**Definition 2 (Minimum Norm Estimator):** Given data $X \in \mathbb{H}^n$ and $\boldsymbol{y} \in \mathbb{R}^n$, the minimum norm estimator $\hat{\theta}$ solves the optimization problem

$$\min_{\theta \in \mathbb{H}} \quad \|\theta\|^2$$

such that $\quad \|X\theta - \boldsymbol{y}\|^2 = \min_{\beta} \|X\beta - \boldsymbol{y}\|^2.$

# 3. Definitions and Notations

●Organize the variables introduced above:

H：随机变量$x$所在的无限维线性空间

$x$：H中的随机变量

$\Sigma$：协方差矩阵$\mathrm{E}[xx']$

$\theta^*$：最优估计值（$y = x'\theta^* + \varepsilon$）

$\varepsilon$：噪声（$\mathrm{E}\varepsilon = 0$）

$\hat{\theta}$：最小范数估计值

# 3. Definitions and Notations

We use $\mu_1(\Sigma) \geq \mu_2(\Sigma) \geq \cdots$ to denote the eigenvalues of $\Sigma$ in descending order, and we denote the operator norm of $\Sigma$ by $\|\Sigma\|$. We use $I$ to denote the identity operator on $\mathbb{H}$ and $I_n$ to denote the $n \times n$ identity matrix.

**Definition 3 (Effective Ranks):** For the covariance operator $\Sigma$, define $\lambda_i = \mu_i(\Sigma)$ for $i = 1, 2, \ldots$. If $\sum_{i=1}^{\infty} \lambda_i < \infty$ and $\lambda_{k+1} > 0$ for $k \geq 0$, define

$$r_k(\Sigma) = \frac{\sum_{i>k} \lambda_i}{\lambda_{k+1}}, \qquad R_k(\Sigma) = \frac{\left(\sum_{i>k} \lambda_i\right)^2}{\sum_{i>k} \lambda_i^2}.$$

# 4.Theorems and Conclusions

**Theorem 1.** *For any* $\sigma_x$*, there are* $b, c, c_1 > 1$ *for which the following holds. Consider a linear regression problem from Definition 1. Define*

$$k^* = \min\{k \geq 0 : r_k(\Sigma) \geq bn\},$$

*where the minimum of the empty set is defined as* $\infty$*. Suppose that* $\delta < 1$ *with* $\log(1/\delta) < n/c$*. If* $k^* \geq n/c_1$*, then* $\mathbb{E}R(\hat{\theta}) \geq \sigma^2/c$*. Otherwise,*

$$R(\hat{\theta}) \leq c\left(\|\theta^*\|^2\|\Sigma\|\max\left\{\sqrt{\frac{r_0(\Sigma)}{n}}, \frac{r_0(\Sigma)}{n}, \sqrt{\frac{\log(1/\delta)}{n}}\right\}\right.$$
$$\left. + c\log(1/\delta)\sigma_y^2\left(\frac{k^*}{n} + \frac{n}{R_{k^*}(\Sigma)}\right)\right)$$

*with probability at least* $1 - \delta$*, and*

$$\mathbb{E}R(\hat{\theta}) \geq \frac{\sigma^2}{c}\left(\frac{k^*}{n} + \frac{n}{R_{k^*}(\Sigma)}\right).$$

*Moreover, there are universal constants* $a_1, a_2, n_0$ *such that, for all* $n \geq n_0$*, for all* $\Sigma$*, and for all* $t \geq 0$*, there is a* $\theta^*$ *with* $\|\theta^*\| = t$ *such that, for* $x \sim \mathcal{N}(0, \Sigma)$ *and* $y|x \sim \mathcal{N}(x^\top\theta^*, \|\theta^*\|^2\|\Sigma\|)$ *with probability at least* $1/4$*,*

$$R(\hat{\theta}) \geq \frac{1}{a_1}\|\theta^*\|^2\|\Sigma\|\mathbb{1}\left[\frac{r_0(\Sigma)}{n\log(1 + r_0(\Sigma))} \geq a_2\right].$$

# 4.Theorems and Conclusions

- **From Theorem 1 we know:**

- $r_0(\Sigma)$ should be small compared with the sample size $n$ (from the first term)

- $r_{k*}(\Sigma)$ and $R_{k*}(\Sigma)$ should be large compared with $n$.

- the **number** of nonzero eigenvalues should be **large** compared with $n$

- they should have a **small sum** compared with $n$

- there should be **many** eigenvalues **no larger** than $\lambda_{k*}$

- If the **number** of these small eigenvalues is **not much larger** than $n$, then they should be roughly **equal**

- they can be more **asymmetric(非对称)** if there are many more of them.

# 4.Theorems and Conclusions

**Theorem 2.** **(Two Examples)**

1) *If* $\mu_k(\Sigma) = k^{-\alpha} \ln^{-\beta}((k+1)e/2)$, *then $\Sigma$ is benign if and only if $\alpha = 1$ and $\beta > 1$.*

Theorem 2.1 shows that, for infinite-dimensional data with a fixed covariance, benign overfitting occurs if and only if the eigenvalues of the covariance operator decay just slowly enough for their sum to remain finite.

# 4.Theorems and Conclusions

Since rescaling $X$ affects the accuracy of the least norm interpolant in an obvious way, we may assume without loss of generality that $\|\Sigma\| = 1$. If we restrict our attention to this case, then informally, Theorem 1 implies that, when the covariance operator for data with $n$ examples is $\Sigma_n$, the least norm interpolant converges if $\frac{r_0(\Sigma_n)}{n} \to 0$, $\frac{k_n^*}{n} \to 0$, and $\frac{n}{R_{k_n^*}(\Sigma_n)} \to 0$ and only if $\frac{r_0(\Sigma_n)}{n\log(1+r_0(\Sigma_n))} \to 0$, $\frac{k_n^*}{n} \to 0$, and $\frac{n}{R_{k_n^*}(\Sigma_n)} \to 0$, where $k_n^* = \min\{k \geq 0 : r_k(\Sigma_n) \geq bn\}$ for the universal constant $b$ in Theorem 1.

**Definition 4:** A sequence of covariance operators $\Sigma_n$ with $\|\Sigma_n\| = 1$ is benign if

$$\lim_{n\to\infty} \frac{r_0(\Sigma_n)}{n} = \lim_{n\to\infty} \frac{k_n^*}{n} = \lim_{n\to\infty} \frac{n}{R_{k_n^*}(\Sigma_n)} = 0.$$

$$R(\hat{\theta}) \leq c \left( \|\theta^*\|^2 \|\Sigma\| \max\left\{ \sqrt{\frac{r_0(\Sigma)}{n}}, \frac{r_0(\Sigma)}{n}, \sqrt{\frac{\log(1/\delta)}{n}} \right\} \right.$$
$$\left. + c\log(1/\delta)\sigma_y^2 \left( \frac{k^*}{n} + \frac{n}{R_{k^*}(\Sigma)} \right) \right)$$

$$\mathbb{E}R(\hat{\theta}) \geq \frac{\sigma^2}{c} \left( \frac{k^*}{n} + \frac{n}{R_{k^*}(\Sigma)} \right).$$

$$R(\hat{\theta}) \geq \frac{1}{a_1} \|\theta^*\|^2 \|\Sigma\| \mathbb{1}\left[ \frac{r_0(\Sigma)}{n\log(1+r_0(\Sigma))} \geq a_2 \right].$$

# 4.Theorems and Conclusions

2) If

$$\mu_k(\Sigma_n) = \begin{cases} \gamma_k + \epsilon_n & \text{if } k \le p_n, \\ 0 & \text{otherwise} \end{cases}$$

*and $\gamma_k = \Theta(\exp(-k/\tau))$, then $\Sigma_n$ with $\|\Sigma_n\| = 1$ is benign if and only if $p_n = \omega(n)$ and $ne^{-o(n)} = \epsilon_n p_n = o(n)$. Furthermore, for $p_n = \Omega(n)$ and $\epsilon_n p_n = ne^{-o(n)}$,*

$$R(\hat{\theta}) = O\left( \frac{\epsilon_n p_n + 1}{n} + \frac{\ln(n/(\epsilon_n p_n))}{n} + \max\left\{ \frac{1}{n}, \frac{n}{p_n} \right\} \right).$$

Theorem 2.2 shows that the situation is very different if the data have finite dimension and a small amount of isotropic noise is added to the covariates. In that case, even if the eigenvalues of the original covariance operator (before the addition of isotropic noise) decay very rapidly, benign overfitting occurs if and only if both the dimension is large compared with the sample size and the isotropic component of the covariance is sufficiently small—but not exponentially small—compared with the sample size.

# 4. Theorems and Conclusions

- **Tension:**

  - the **slow decay** of eigenvalues that is needed for $k/n + n/R_k$ to be small

  - the **summability** of eigenvalues that is needed for $r_0(\Sigma)/n$ to be small

# 4.Theorems and Conclusions

- **Explain**

- 1）In the **infinite**-dimensional setting, **slow** decay of the eigenvalues suffices—decay just **fast enough** to ensure **summability**—as shown by Theorem 2.1.

- (Example)

**Theorem 31.** *Define* $\lambda_{k,n} := \mu_k(\Sigma_n)$ *for all* $k, n$.

1. *If* $\lambda_{k,n} = k^{-\alpha} \ln^{-\beta}(k+1)$, *then* $\Sigma_n$ *is benign iff* $\alpha = 1$ *and* $\beta > 1$.

2. *If* $\lambda_{k,n} = k^{-(1+\alpha_n)}$, *then* $\Sigma_n$ *is benign iff* $\omega(1/n) = \alpha_n = o(1)$. *Furthermore,*

$$R(\hat{\theta}) = \Theta\left(\min\left\{\frac{1}{\alpha_n n} + \alpha_n, 1\right\}\right).$$

# 4.Theorems and Conclusions

- **Explain**

- 2）Consider a **finite**-dimensional setting (which ensures that the eigenvalues are **summable**), and in this case, arbitrarily slow decay is possible.

- (Example)

3. If

$$\lambda_{k,n} = \begin{cases} k^{-\alpha} & \text{if } k \leq p_n, \\ 0 & \text{otherwise,} \end{cases}$$

then $\Sigma_n$ is benign iff either $0 < \alpha < 1$, $p_n = \omega(n)$ and $p_n = o\left(n^{1/(1-\alpha)}\right)$ or $\alpha = 1$, $p_n = e^{\omega(\sqrt{n})}$ and $p_n = e^{o(n)}$.

4. If

$$\lambda_{k,n} = \begin{cases} \gamma_k + \epsilon_n & \text{if } k \leq p_n, \\ 0 & \text{otherwise,} \end{cases}$$

and $\gamma_k = \Theta(\exp(-k/\tau))$, then $\Sigma_n$ is benign iff $p_n = \omega(n)$ and $ne^{-o(n)} = \epsilon_n p_n = o(n)$. Furthermore, for $p_n = \Omega(n)$ and $\epsilon_n p_n = ne^{-o(n)}$,

$$R(\hat{\theta}) = O\left(\frac{\epsilon_n p_n + 1}{n} + \frac{\ln(n/(\epsilon_n p_n))}{n} + \max\left\{\frac{1}{n}, \frac{n}{p_n}\right\}\right).$$

# 5.Proof

**●5.1.Headline:**

●   a standard decomposition of the excessrisk into two pieces:

●    1) a term that corresponds to the distortion that is introduced by viewing $\theta^{*}$ through the lens of **the finite sample**

●    2) a term that corresponds to the distortion introduced by the **noise**

$$\varepsilon = y - X\hat{\theta}$$

# 5.Proof

● **5.2.可能用到的数学知识**

● 广义逆矩阵

**定理 1** 设 $A$ 是数域 $K$ 上 $s \times n$ 非零矩阵,则矩阵方程

$$AXA = A \tag{2}$$

一定有解。如果 $\operatorname{rank}(A) = r$,并且

$$A = P \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix} Q, \tag{3}$$

其中 $P, Q$ 分别是 $K$ 上 $s$ 级、$n$ 级可逆矩阵,那么矩阵方程(2)的通解为

$$X = Q^{-1} \begin{pmatrix} I_r & B \\ C & D \end{pmatrix} P^{-1}. \tag{4}$$

其中 $B, C, D$ 分别是数域 $K$ 上任意的 $r \times (s-r), (n-r) \times r, (n-r) \times (s-r)$ 矩阵。

# 5.Proof

**定义 1**　设 $A$ 是数域 $K$ 上 $s \times n$ 矩阵,矩阵方程 $AXA = A$ 的每一个解都称为 $A$ 的一个广义逆矩阵,简称为 $A$ 的广义逆,用 $A^-$ 表示 $A$ 的任意一个广义逆。

从定义 1 得出

$$AA^-A = A. \tag{10}$$

从定理 1 得出,当 $A \neq 0$ 时,设 $\text{rank}(A) = r$,且

$$A = P\begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix}Q.$$

则

$$A^- = Q^{-1}\begin{pmatrix} I_r & B \\ C & D \end{pmatrix}P^{-1}. \tag{11}$$

从定义 1 得出,任意一个 $n \times s$ 矩阵都是 $0_{s \times n}$ 的广义逆。

**定理 2**　(非齐次线性方程组的相容性定理) 非齐次线性方程组 $AX = \beta$ 有解的充分必要条件是

$$\beta = AA^-\beta. \tag{12}$$

**证明**　必要性。设 $AX = \beta$ 有解 $\alpha$,则

$$\beta = A\alpha = AA^-A\alpha = AA^-\beta.$$

充分性。设 $\beta = AA^-\beta$,则 $A^-\beta$ 是 $AX = \beta$ 的解。　∎

**定理 3**　(非齐次线性方程组的解的结构定理)非齐次线性方程组 $AX = \beta$ 有解时,它的通解为

$$X = A^-\beta. \tag{13}$$

**定理 4**　(齐次线性方程组的解的结构定理)数域 $K$ 上 $n$ 元齐次线性方程组 $AX = 0$ 的通解为

$$X = (I_n - A^-A)Z, \tag{20}$$

其中 $A^-$ 是 $A$ 的任意给定的一个广义逆,$Z$ 取遍 $K^n$ 中任意列向量。

**证明**　任取 $Z \in K^n$,有

$$A[(I_n - A^-A)Z] = (A - AA^-A)Z = (A - A)Z = 0,$$

因此 $X = (I_n - A^-A)Z$ 是齐次线性方程组 $AX = 0$ 的解。

反之,设 $\eta$ 是 $AX = 0$ 的一个解,则

$$(I_n - A^-A)\eta = \eta - A^-A\eta = \eta.$$

综上所述,$X = (I_n - A^-A)Z$ 是齐次线性方程组 $AX = 0$ 的通解。　∎

**推论 1**　设数域 $K$ 上 $n$ 元非齐次线性方程组 $AX = \beta$ 有解,则它的通解为

$$X = A^-\beta + (I_n - A^-A)Z, \tag{21}$$

其中 $A^-$ 是 $A$ 的任意给定的一个广义逆,$Z$ 取遍 $K^n$ 中任意列向量。

**证明**　由于 $A^-\beta$ 是 $AX = \beta$ 的一个解,且 $(I_n - A^-A)Z$ 是导出方程组 $AX = 0$ 的通解,因此 $X = A^-\beta + (I_n - A^-A)Z$ 是 $AX = \beta$ 的通解。　∎

# 5.Proof

**定义 2** 设 $A$ 是复数域上 $s \times n$ 矩阵,矩阵方程组

$$\begin{cases} AXA = A, \\ XAX = X, \\ (AX)^* = AX, \\ (XA)^* = XA, \end{cases} \tag{22}$$

称为 $A$ 的 **Penrose 方程组**,它的解称为 $A$ 的 **Moore-Penrose 广义逆**,记作 $A^+$。(22)式中 $(AX)^*$ 表示把 $AX$ 的每个元素取共轭复数得到的矩阵再转置。

**定理 5** 如果 $A$ 是复数域上 $s \times n$ 非零矩阵,$A$ 的 Penrose 方程组总是有解,并且它的解唯一。设 $A = BC$,其中 $B$、$C$ 分别是列满秩与行满秩矩阵,则 Penrose 方程组的唯一解是

$$X = C^* (CC^*)^{-1} (B^* B)^{-1} B^*. \tag{23}$$

**证明** 把(23)式代入 Penrose 方程组的每一个方程,验证每一个方程都变成恒等式:

$$AXA = (BC)C^*(CC^*)^{-1}(B^*B)^{-1}B^*(BC) = BC = A,$$

$$XAX = C^*(CC^*)^{-1}(B^*B)^{-1}B^*(BC)C^*(CC^*)^{-1}(B^*B)^{-1}B^*$$
$$= C^*(CC^*)^{-1}(B^*B)^{-1}B^* = X,$$

$$(AX)^* = X^*A^* = B(B^*B)^{-1}(CC^*)^{-1}CC^*B^*$$
$$= B(B^*B)^{-1}B^* = B(CC^*)(CC^*)^{-1}(B^*B)^{-1}B^* = AX,$$

$$(XA)^* = A^*X^* = C^*B^*B(B^*B)^{-1}(CC^*)^{-1}C$$
$$= C^*(CC^*)^{-1}C = C^*(CC^*)^{-1}(B^*B)^{-1}(B^*B)C = XA.$$

因此(12)式的确是 Penrose 方程组的解。

下面证解的唯一性。设 $X_1$ 和 $X_2$ 都是 Penrose 方程组的解。则

$$X_1 = X_1 A X_1 = X_1 (A X_2 A) X_1 = X_1 (A X_2)(A X_1)$$

$$= X_1 (A X_2)^* (A X_1)^* = X_1 (A X_1 A X_2)^* = X_1 X_2^* (A X_1 A)^*$$
$$= X_1 X_2^* A^* = X_1 (A X_2)^* = X_1 A X_2 = X_1 (A X_2 A) X_2$$
$$= (X_1 A)(X_2 A) X_2 = (X_1 A)^* (X_2 A)^* X_2 = (X_2 A X_1 A)^* X_2$$
$$= (X_2 A)^* X_2 = X_2 A X_2 = X_2.$$

# 5.Proof

## ●5.3. Lemma

**Lemma 1.** $r_k(\Sigma) \geq 1$, $r_k^2(\Sigma) = r_k(\Sigma^2) R_k(\Sigma)$, and

$$r_k(\Sigma^2) \leq r_k(\Sigma) \leq R_k(\Sigma) \leq r_k^2(\Sigma).$$

Notice that $r_0(I_p) = R_0(I_p) = p$. More generally, if all of the nonzero eigenvalues of $\Sigma$ are identical, then $r_0(\Sigma) = R_0(\Sigma) = \text{rank}(\Sigma)$. For $\Sigma$ with finite rank, we can express both $r_0(\Sigma)$ and $R_0(\Sigma)$ as a product of the rank and a notion of symmetry. In particular, for $\text{rank}(\Sigma) = p$, we can write

$$r_0(\Sigma) = \text{rank}(\Sigma) s(\Sigma), \qquad R_0(\Sigma) = \text{rank}(\Sigma) S(\Sigma),$$

$$\text{with } s(\Sigma) = \frac{\frac{1}{p}\sum_{i=1}^{p} \lambda_i}{\lambda_1}, \qquad S(\Sigma) = \frac{\left(\frac{1}{p}\sum_{i=1}^{p} \lambda_i\right)^2}{\frac{1}{p}\sum_{i=1}^{p} \lambda_i^2}.$$

Both notions of symmetry $s$ and $S$ lie between $1/p$ (when $\lambda_2 \to$ 0) and 1 (when the $\lambda_i$ are all equal).

# 5.Proof

**Lemma 2.** *The excess risk of the minimum norm estimator satisfies* $R(\hat{\theta}) \le 2\theta^{*\top} B\theta^* + c\sigma^2 \log(1/\delta)\,\mathrm{tr}(C)$ *with probability at least* $1-\delta$ *over* $\epsilon$, *and* $\mathbb{E}_\varepsilon R(\hat{\theta}) \ge \theta^{*\top} B\theta^* + \sigma^2\,\mathrm{tr}(C)$, *where*

$$B = \left(I - X^\top \left(XX^\top\right)^{-1}X\right)\Sigma\left(I - X^\top \left(XX^\top\right)^{-1}X\right),$$

$$C = \left(XX^\top\right)^{-1}X\Sigma X^\top \left(XX^\top\right)^{-1}.$$

Proof

$\cdots$

Le2 $R(\hat{\theta}) = \mathbb{E}_{x,y}(y - x'\hat{\theta})^2 - \mathbb{E}(y - x'\theta^*)^2$

$= \mathbb{E}_{x,y}\left(y - x'\theta^* + x'(\theta^* - \hat{\theta})\right)^2 - \mathbb{E}(y - x'\theta^*)^2$

$= \mathbb{E}_x\left(x'(\theta^* - \hat{\theta})\right)^2$

$y = x\theta^* + \varepsilon = x\hat{\theta}$ 得到

$\hat{\theta} = x'(xx')^{-1}x\theta^* + x'(xx')^{-1}\varepsilon$

$\mathbb{E}_x\left(x'(\theta^* - \hat{\theta})\right)^2 = \mathbb{E}_x\left(x'\left((I - x'(xx')^{-1}x)\theta^* - x'(xx')^{-1}\varepsilon\right)\right)^2$

$\le 2\mathbb{E}_x\left(x'(I - x'(xx')^{-1}x)\theta^*\right)^2 + 2\mathbb{E}_x\left(x'x'(xx')^{-1}\varepsilon\right)^2$

$= 2\mathbb{E}_x\left(\theta^{*'}(I - x'(xx')^{-1}x)\underbrace{xx'}_{\Sigma}(I - x'(xx')^{-1}x)\theta^*\right)$

$\quad + 2\mathbb{E}_x\left(\varepsilon'(xx')^{-1}x\underbrace{xx'}_{\Sigma}x'(xx')^{-1}\varepsilon\right)$

$= 2\theta^{*'}B\theta^* + 2\varepsilon'C\varepsilon$

$\varepsilon$ 是 $\delta^2$ 次高斯的, 即 $\forall \lambda \in R$

$$\mathbb{E}[e^{\lambda\varepsilon}|x] \le e^{\frac{\lambda^2\delta^2}{2}}$$

以及 $1 - e^{-t}$ 的不等式

$\varepsilon'C\varepsilon \le \delta^2\,\mathrm{tr}(C) + 2\delta^2 + |c| + 2\delta^2\sqrt{|c|t^2 + \mathrm{tr}(c^2)t}$ (1)

# 5.Proof

**Lemma 2.** *The excess risk of the minimum norm estimator satisfies* $R(\hat{\theta}) \leq 2\theta^{*\top} B\theta^* + c\sigma^2 \log(1/\delta)\operatorname{tr}(C)$ *with probability at least* $1-\delta$ *over* $\epsilon$, *and* $\mathbb{E}_\epsilon R(\hat{\theta}) \geq \theta^{*\top} B\theta^* + \sigma^2\operatorname{tr}(C)$, *where*

$$B = \left(I - X^\top \left(XX^\top\right)^{-1}X\right)\Sigma\left(I - X^\top\left(XX^\top\right)^{-1}X\right),$$
$$C = \left(XX^\top\right)^{-1}X\Sigma X^\top\left(XX^\top\right)^{-1}.$$

$\forall v \in k^n,$

$$\frac{|C\alpha|}{|\alpha|} = \sqrt{\frac{(C\alpha, C\alpha)}{(\alpha, \alpha)}} = \sqrt{\frac{\alpha' C' C \alpha}{\alpha' \alpha}}$$

$C$ 是对称矩阵，同此，存在正交矩阵 $P$ s.t.

$$P C P' = diag\{\mu_1, \mu_2 \cdots \mu_n\} := D$$

$$\frac{|C\alpha|}{|\alpha|} = \sqrt{\frac{\alpha' P' D^2 P\alpha}{\alpha'\alpha}} = \sqrt{\frac{(\mu_1\theta_1)^2 + (\mu_2\theta_2)^2 + \cdots + (\mu_n\theta_n)^2}{\alpha_1^2 + \alpha_2^2 + \cdots + \alpha_n^2}}$$

$$\leq \sqrt{\frac{\operatorname{tr}(C^2)(\theta_1^2 + \cdots + \theta_n^2)}{\alpha_1^2 + \cdots + \alpha_n^2}} = \sqrt{\operatorname{tr}(C^2)} \cdot \frac{|P\alpha|}{|\alpha|} = \sqrt{\operatorname{tr}(C^2)}$$

$$\operatorname{tr}(C^2) = \mu_1^2 + \mu_2^2 + \cdots + \mu_n^2 \leq (\mu_1 + \mu_2 + \cdots + \mu_n)^2$$

$= (\operatorname{tr} C)^2$, 因此

$$\frac{|C\alpha|}{|\alpha|} \leq \operatorname{tr} C. \quad 结合 |C| = \sup_{\alpha \in k^n} \frac{|C\alpha|}{|\alpha|}, \text{ 因此}$$

$$|C| \leq \operatorname{tr} C$$

代入 (1) 式有，以大于 $1-e^{-t}$ 的概率

$$\varepsilon' C \varepsilon \leq \beta^2 \operatorname{tr}(C)(2t+1) + 2\beta^2 \operatorname{tr}(C)\sqrt{t^2 + t}$$

$$\leq (4t+2)\beta^2 \operatorname{tr}(C)$$

取 $t = \log(\frac{1}{\delta})$，有，以大于 $1-\delta$ 的概率

$$R(\hat{\theta}) \leq 2\varepsilon' B\varepsilon + 2\varepsilon' C\varepsilon$$

$$\leq 2\varepsilon' B\varepsilon + 4\left(2\log\frac{1}{\delta} + 1\right)\beta^2 \operatorname{tr}(C)$$

# 5.Proof

**Lemma 2.** *The excess risk of the minimum norm estimator satisfies* $R(\hat{\theta}) \le 2\theta^{*\top} B \theta^* + c\sigma^2 \log(1/\delta) \operatorname{tr}(C)$ *with probability at least* $1 - \delta$ *over* $\epsilon$, *and* $\mathbb{E}_\epsilon R(\hat{\theta}) \ge \theta^{*\top} B \theta^* + \sigma^2 \operatorname{tr}(C)$, *where*

$$B = \left(I - X^\top \left(XX^\top\right)^{-1} X\right) \Sigma \left(I - X^\top \left(XX^\top\right)^{-1} X\right),$$

$$C = \left(XX^\top\right)^{-1} X \Sigma X^\top \left(XX^\top\right)^{-1}.$$

$$\mathbb{E}R(\hat{\theta}) = \mathbb{E}_{x,\varepsilon}\left[(x'B\theta^*)^2 + (x'X'(XX')^{-1}\varepsilon)^2\right] \mid 交叉项是关于\varepsilon 的线性函数均值为0)$$

$$= \mathbb{E}_{x,\varepsilon}\, \theta^{*'} B' x x' B \theta^* + \mathbb{E}_{x,\varepsilon}\, \varepsilon'(XX')^{-1} X\, x x'\, X'(XX')^{-1}\varepsilon$$

$$= \theta^{*'}\, B\, \theta^* + \mathbb{E}_\varepsilon\, \varepsilon'(XX')^{-1} X \Sigma x'(XX')^{-1}\varepsilon$$

$$= \theta^{*'}\, B\, \theta^* + \operatorname{tr}\left((XX')^{-1} X \Sigma X'(XX')^{-1} \mathbb{E}[\varepsilon\varepsilon' \mid x]\right)$$

$$= \theta^{*'} B \theta^* + \operatorname{tr}\left(C\, \mathbb{E}[\varepsilon\varepsilon' \mid x]\right) \le \theta^{*'} B \theta^* + \delta^2 \operatorname{tr}(C)$$

# 5.Proof

**Lemma 3.** *Consider a covariance operator $\Sigma$ with $\lambda_i = \mu_i(\Sigma)$ and $\lambda_n > 0$. Write its spectral decomposition $\Sigma = \sum_j \lambda_j v_j v_j^\top$, where the orthonormal $v_j \in \mathbb{H}$ are the eigenvectors corresponding to the $\lambda_j$. For $i$ with $\lambda_i > 0$, define $z_i = X v_i / \sqrt{\lambda_i}$. Then,*

$$\operatorname{tr}(C) = \sum_i \left[ \lambda_i^2 z_i^\top \left( \sum_j \lambda_j z_j z_j^\top \right)^{-2} z_i \right],$$

*and these $z_i \in \mathbb{R}^n$ are independent $\sigma_x^2$ sub-Gaussian. Furthermore, for any $i$ with $\lambda_i > 0$, we have*

$$\lambda_i^2 z_i^\top \left( \sum_j \lambda_j z_j z_j^\top \right)^{-2} z_i = \frac{\lambda_i^2 z_i^\top A_{-i}^{-2} z_i}{(1 + \lambda_i z_i^\top A_{-i}^{-1} z_i)^2},$$

*where $A_{-i} = \sum_{i \neq i} \lambda_j z_j z_j^\top$.*

Le3 由Def 1中前的Assumption 2 可知，$\frac{\alpha' v_i}{\sqrt{\lambda_i}}$ 为 $\partial_x^2$ 次高斯变量

注意到 $X v_i = \sqrt{\lambda_i}\, \delta_i$;

$\{v_j\}$ 为 $H$ 上前一组标准正交基。因此，$X$ 的一行可以表示为

$$x_k' = \sum_j (x_k' v_j) v_j'$$

$$(\delta_i)_k = \frac{x_k' v_i}{\sqrt{\lambda_i}}$$

所以 $x_k' = \sum_j \sqrt{\lambda_j} (\delta_j)_k v_j'$

$$X = \begin{pmatrix} \sum_j \sqrt{\lambda_j} (\delta_j)_1 v_j' \\ \cdots \\ \sum_j \sqrt{\lambda_j} (\delta_j)_n v_j' \end{pmatrix} = \sum_j \sqrt{\lambda_j}\, \delta_j v_j'$$

# 5.Proof

**Lemma 3.** *Consider a covariance operator $\Sigma$ with $\lambda_i = \mu_i(\Sigma)$ and $\lambda_n > 0$. Write its spectral decomposition $\Sigma = \sum_j \lambda_j v_j v_j^\top$, where the orthonormal $v_j \in \mathbb{H}$ are the eigenvectors corresponding to the $\lambda_j$. For $i$ with $\lambda_i > 0$, define $z_i = X v_i / \sqrt{\lambda_i}$. Then,*

$$\operatorname{tr}(C) = \sum_i \left[ \lambda_i^2 z_i^\top \left( \sum_j \lambda_j z_j z_j^\top \right)^{-2} z_i \right],$$

*and these $z_i \in \mathbb{R}^n$ are independent $\sigma_x^2$ sub-Gaussian. Furthermore, for any $i$ with $\lambda_i > 0$, we have*

$$\lambda_i^2 z_i^\top \left( \sum_j \lambda_j z_j z_j^\top \right)^{-2} z_i = \frac{\lambda_i^2 z_i^\top A_{-i}^{-2} z_i}{(1 + \lambda_i z_i^\top A_{-i}^{-1} z_i)^2},$$

*where $A_{-i} = \sum_{i \neq i} \lambda_j z_j z_j^\top$.*

$$XX' = \sum_i \sum_j \sqrt{\lambda_i} \sqrt{\lambda_j}\, \delta_j v_j' v_i\, \delta_j'$$

$$= \sum_j \lambda_j \delta_j \delta_j'$$

$$X \Sigma X' = \left( \sum_i \sqrt{\lambda_i}\, \delta_i v_i' \right) \left( \sum_j \lambda_j v_j v_j' \right) \left( \sum_k \sqrt{\lambda_k}\, v_k \delta_k' \right)$$

$$= \sum_i \sum_j \sum_k \sqrt{\lambda_i}\sqrt{\lambda_k}\, \lambda_j\, \delta_i (v_i' v_j)(v_j' v_k)\, \delta_k'$$

$$= \sum_j \lambda_j^2\, \delta_j\, \delta_j'$$

$$\operatorname{tr}(C) = \operatorname{tr}\left( (XX')^{-1} X\Sigma X' (XX')^{-1} \right)$$

$$= \operatorname{tr}\left( (XX')^{-2} X\Sigma X' \right)$$

$$= \operatorname{tr}\left( \left( \sum_j \lambda_j \delta_j \delta_j' \right)^{-2} \left( \sum_j \lambda_j^2 \delta_j \delta_j' \right) \right)$$

$$= \sum_j \lambda_j^2\, \operatorname{tr}\left( \left( \sum_i \lambda_i \delta_i \delta_i' \right)^{-2} \delta_j \delta_j' \right)$$

$$= \sum_j \lambda_j^2\, \operatorname{tr}\left( \delta_j' \left( \sum_i \lambda_i \delta_i \delta_i' \right)^{-2} \delta_j \right)$$

$$= \sum_j \lambda_j^2 \delta_j' \left( \sum_i \lambda_i \delta_i \delta_i' \right)^{-2} \delta_j$$

# 5.Proof

**Lemma 3.** *Consider a covariance operator $\Sigma$ with $\lambda_i = \mu_i(\Sigma)$ and $\lambda_n > 0$. Write its spectral decomposition $\Sigma = \sum_j \lambda_j v_j v_j^\top$, where the orthonormal $v_j \in \mathbb{H}$ are the eigenvectors corresponding to the $\lambda_j$. For $i$ with $\lambda_i > 0$, define $z_i = X v_i / \sqrt{\lambda_i}$. Then,*

$$\operatorname{tr}(C) = \sum_i \left[ \lambda_i^2 z_i^\top \left( \sum_j \lambda_j z_j z_j^\top \right)^{-2} z_i \right],$$

*and these $z_i \in \mathbb{R}^n$ are independent $\sigma_x^2$ sub-Gaussian. Furthermore, for any $i$ with $\lambda_i > 0$, we have*

$$\lambda_i^2 z_i^\top \left( \sum_j \lambda_j z_j z_j^\top \right)^{-2} z_i = \frac{\lambda_i^2 z_i^\top A_{-i}^{-2} z_i}{(1 + \lambda_i z_i^\top A_{-i}^{-1} z_i)^2},$$

*where $A_{-i} = \sum_{i \neq i} \lambda_j z_j z_j^\top$.*

根据 Sherman-Morrison 公式： $(A + \alpha \alpha')^{-1} = A^{-1} - \dfrac{A^{-1} x x' A^{-1}}{1 + \alpha' A^{-1} \alpha}$

$$\left( \sum_j \lambda_j \partial_j \partial_j' \right)^{-1} = \left( A_{-i} + \lambda_i \partial_i \partial_i' \right)^{-1}$$

$$= A_{-i}^{-1} - \frac{\lambda_i A_{-i}^{-1} \partial_i \partial_i' A_{-i}^{-1}}{1 + \lambda_i \partial_i' A_{-i}^{-1} \partial_i}$$

通过计算可得 (过程省去)

$$\lambda_i^2 \partial_i' \left( \sum_j \lambda_j \partial_j \partial_j' \right) \partial_i = \frac{\lambda_i^2 \partial_i' A_{-i}^{-2} \partial_i}{\left( 1 + \lambda_i \partial_i' A_{-i}^{-1} \partial_i \right)^2}$$

整合过后：

$$\operatorname{tr}(c) = \sum_j \frac{\lambda_j^2 \partial_j' A_{-j}^{-2} \partial_j}{\left( 1 + \lambda_j \partial_j' A_{-j}^{-1} \partial_j \right)^2}$$

# 5.Proof

● **Define:**

$$A = \sum_i \lambda_i z_i z_i^\top, \quad A_{-i} = \sum_{j \neq i} \lambda_j z_j z_j^\top, \quad A_k = \sum_{i > k} \lambda_i z_i z_i^\top,$$

● The next step is to show that eigenvalues of $A$, $A_{-i}$, and $A_k$ are **concentrated**.

# 5.Proof

**Lemma 4.** *There is a constant $c$ such that, for any $k \geq 0$ with probability at least $1 - 2e^{-n/c}$,*

$$\frac{1}{c}\sum_{i>k}\lambda_i - c\lambda_{k+1}n \leq \mu_n(A_k) \leq \mu_1(A_k) \leq c\left(\sum_{i>k}\lambda_i + \lambda_{k+1}n\right).$$

**Lemma 25** ($\epsilon$-net argument). *Suppose $A \in \mathbb{R}^{n\times n}$ is a symmetric matrix, and $\mathcal{N}_\epsilon$ is an $\epsilon$-net on the unit sphere $\mathcal{S}^{n-1}$ in the Euclidean norm, where $\epsilon < \frac{1}{2}$. Then*

$$\|A\| \leq (1-\epsilon)^{-2}\max_{x\in\mathcal{N}_\epsilon}|x^\top A x|.$$

*Proof.* Denote the eigenvalues of $A$ as $\lambda_1, \ldots, \lambda_n$ and assume $|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_n|$. Denote the first eigenvector of $A$ as $v \in \mathcal{S}^{n-1}$, and take $\Delta v \in \mathbb{R}^n$ such that $v + \Delta v \in \mathcal{N}_\epsilon$ and $\|\Delta v\| \leq \epsilon$. Denote the coordinates of $\Delta v$ in the eigenbasis of $A$ as $\Delta v_1, \ldots, \Delta v_n$. Now we can write

$$\left|(v + \Delta v)^\top A(v + \Delta v)\right| = \left|\lambda_1 + 2\lambda_1\Delta v_1 + \sum_{i=1}^{n}\lambda_i\Delta v_i^2\right|$$

$$= |\lambda_1| \cdot \left|1 + 2\Delta v_1 + \Delta v_1^2 + \sum_{i=2}^{n}\frac{\lambda_i}{\lambda_1}\Delta v_i^2\right|$$

$$\geq |\lambda_1| \cdot \left|1 + 2\Delta v_1 + \Delta v_1^2 - \sum_{i=2}^{n}\Delta v_i^2\right|$$

$$= |\lambda_1| \cdot \left|1 + 2\Delta v_1 + \Delta v_1^2 - \|\Delta v\|^2 + \Delta v_1^2\right|$$

$$= |\lambda_1| \cdot \left|1 + 2\left(\Delta v_1 + \Delta v_1^2\right) - \|\Delta v\|^2\right|$$

$$\geq |\lambda_1| \cdot \left|1 + 2\left(-\|\Delta v\| + (-\|\Delta v\|)^2\right) - \|\Delta v\|^2\right|$$

$$= |\lambda_1| \cdot \left|1 - 2\|\Delta v\| + \|\Delta v\|^2\right|$$

$$\geq |\lambda_1| \cdot \left|1 - 2\epsilon + \epsilon^2\right|$$

$$= \|A\|(1-\epsilon)^2,$$

where the first inequality holds because the $\lambda_i$s are decreasing in magnitude, and the last two inequalities hold since the functions $x + x^2$ and $2x + x^2$ are both increasing on $(-\frac{1}{2}, \infty)$ and $\Delta v_1 \geq -\|\Delta v\| \geq -\epsilon \geq -\frac{1}{2}$. $\qquad\square$

# 5.Proof

**Lemma 4.** *There is a constant $c$ such that, for any $k \geq 0$ with probability at least $1 - 2e^{-n/c}$,*

$$\frac{1}{c}\sum_{i>k}\lambda_i - c\lambda_{k+1}n \leq \mu_n(A_k) \leq \mu_1(A_k) \leq c\left(\sum_{i>k}\lambda_i + \lambda_{k+1}n\right).$$

Le 4

对于 $u \in R^n$ ($\|u\|=1$), $u'Au = \sum \lambda_i (u'\delta_i)^2$, 存在 $C_2>0$, 以至少 $1-2e^{-t}$ 的概率

$$\left| u'Au - \sum \lambda_i \right| \leq C_2 \partial_x^2 \max\left(\lambda_1 t, \sqrt{t\sum \lambda_i^2}\right)$$

设 $N$ 为 $S^{n-1}$ 的 $\frac{1}{4}$-网且 $|N| \leq 9^n$, 以至少 $1-2e^{-t}$ 的概率

$$\left| u'Au - \sum \lambda_i \right| \leq C_2 \partial_x^2 \max\left\{\lambda_1^{(t+n\ln 9)}, \sqrt{(t+n\ln 9)\sum \lambda_i^2}\right\}$$

设 $\hat{v} = \lambda_1(t+n\ln 9) + \sqrt{(t+n\ln 9)\sum \lambda_i^2}$. 存在 $C_3'>0$ s.t.

$$\left| A - \left(\sum_i \lambda_i\right) I_n \right| \leq C_3' \max_{\alpha \in N} \left| \alpha'(A-(\sum_i \lambda_i)I_n)\alpha \right| \leq C_3' \left| u'Au - \sum \lambda_i \right|$$

$$\leq C_3 \partial_x^2 \hat{v}$$

以至少 $1-2e^{-t}$ 概率成立.

$\exists C_4, C_5$ s.t. $t + n\ln 9 \leq C_5 n$ as $t \leq \frac{n}{C_4}$, 此时有

$$\hat{v} \leq C_5\left(\lambda_1 n + \sqrt{n\sum_i \lambda_i^2}\right)$$

由 AM-GM 不等式

$$\sqrt{\lambda_1 n \sum \lambda_i} \leq \frac{1}{2\delta}\lambda_1 n + \frac{\delta}{2}\sum \lambda_i (\forall \delta >0)$$

取 $\delta = C_3 \partial_x^2$, 适当整合不等

$$\hat{v} \leq C_5 \partial_x^2 \lambda_1 n + \frac{1}{2C_3 \partial_x^2}\sum \lambda_i$$

$$\left| A - \left(\sum_i \lambda_i\right)I_n \right| = \max\left\{\left|\mu_1(A)-\sum \lambda_i\right|, \left|\mu_n(A)-\sum \lambda_i\right|\right\}$$

$$\leq \frac{1}{2}\sum \lambda_i + C_6 C_3 \partial_x^4 \lambda_1 n$$

而 $\mu_1(A) - \sum \lambda_i \leq \frac{1}{2}\sum \lambda_i + C_6 C_3 \partial_x^4 \lambda_1 n$

$$\sum \lambda_i - \mu_n(A) \leq \frac{1}{2}\sum \lambda_i + C_6 C_3 \partial_x^4 \lambda_1 n$$

适当调整 $C_1, C_2, \cdots, C_6$, 得到 Lemma 4

# 5.Proof

**Lemma 5** *There are constants $b, c \geq 1$ such that for any $k \geq 0$, with probability at least $1 - 2e^{-n/c}$,*
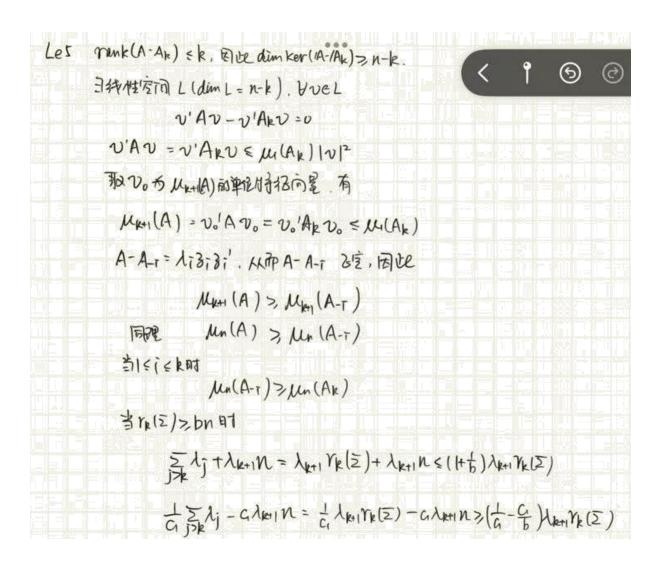
*1. for all $i \geq 1$,*

$$\mu_{k+1}(A_{-i}) \leq \mu_{k+1}(A) \leq \mu_1(A_k) \leq c \left( \sum_{j>k} \lambda_j + \lambda_{k+1} n \right),$$

*2. for all $1 \leq i \leq k$,*

$$\mu_n(A) \geq \mu_n(A_{-i}) \geq \mu_n(A_k) \geq \frac{1}{c} \sum_{j>k} \lambda_j - c\lambda_{k+1} n,$$

*3. if $r_k(\Sigma) \geq bn$, then*

$$\frac{1}{c} \lambda_{k+1} r_k(\Sigma) \leq \mu_n(A_k) \leq \mu_1(A_k) \leq c\lambda_{k+1} r_k(\Sigma).$$

# 5.Proof

Le5  $\text{rank}(A-A_k) \leq k$, 因此 $\dim\ker((A-A_k)) \geq n-k$.

∃线性空间 $L$ $(\dim L = n-k)$, $\forall v \in L$

$$v'Av - v'A_kv = 0$$

$$v'Av = v'A_kv \leq \mu_1(A_k)|v|^2$$

取 $v_0$ 为 $\mu_{k+1}(A)$ 前单位特征向量, 有

$$\mu_{k+1}(A) = v_0'Av_0 = v_0'A_kv_0 \leq \mu_1(A_k)$$

$A - A_T = \lambda_i \delta_i \delta_i'$. 从而 $A - A_T$ 半定, 因此

$$\mu_{k+1}(A) \geq \mu_{k+1}(A_T)$$

同理    $\mu_n(A) \geq \mu_n(A_T)$

当 $1 \leq i \leq k$ 时

$$\mu_n(A_T) \geq \mu_n(A_k)$$

当 $r_k(\bar{\Sigma}) \geq bn$ 时

$$\sum_{j>k} \lambda_j + \lambda_{k+1}n = \lambda_{k+1}r_k(\bar{\Sigma}) + \lambda_{k+1}n \leq (1+\frac{1}{b})\lambda_{k+1}r_k(\bar{\Sigma})$$

$$\frac{1}{c_1}\sum_{j>k}\lambda_j - c_1\lambda_{k+1}n = \frac{1}{c_1}\lambda_{k+1}r_k(\bar{\Sigma}) - c_1\lambda_{k+1}n \geq (\frac{1}{c_1} - \frac{c_1}{b})\lambda_{k+1}r_k(\bar{\Sigma})$$

# 5.Proof

**Lemma** 6  *Suppose $\{\lambda_i\}_i^\infty$ is a non-increasing sequence of non-negative numbers such that $\sum_{i=1}^\infty \lambda_i < \infty$, and $\{\xi_i\}_{i=1}^\infty$ are independent centered $\sigma$-subexponential random variables. Then for some universal constant $a$ for any $t > 0$ with probability at least $1 - 2e^{-t}$*

$$\left| \sum_i \lambda_i \xi_i \right| \leq a\sigma \max\left( t\lambda_1, \sqrt{t \sum_i \lambda_i^2} \right).$$

**Corollary 1.** *Suppose that $z \in \mathbb{R}^n$ is a centered random vector with independent $\sigma^2$ sub-Gaussian coordinates with unit variances, $\mathscr{L}$ is a random subspace of $\mathbb{R}^n$ of codimension $k$, and $\mathscr{L}$ is independent of $z$. Then, for some universal constant $a$ and any $t > 0$, with probability at least $1 - 3e^{-t}$,*

$$\|z\|^2 \leq n + a\sigma^2(t + \sqrt{nt}),$$
$$\|\Pi_{\mathscr{L}} z\|^2 \geq n - a\sigma^2(k + t + \sqrt{nt}),$$

*where $\Pi_{\mathscr{L}}$ is the orthogonal projection on $\mathscr{L}$.*

# 5.Proof

**Upper Bound on the Trace Term. Lemma** [7] *There are constants* $b, c \geq 1$ *such that, if* $0 \leq k \leq n/c$, $r_k(\Sigma) \geq bn$, *and* $l \leq k$, *then with probability at least* $1 - 7e^{-n/c}$,

$$\mathrm{tr}(C) \leq c\left(\frac{l}{n} + n\frac{\sum_{i>l}\lambda_i^2}{(\sum_{i>k}\lambda_i)^2}\right).$$

# 5.Proof

**Lower bound on tr(C)**

**Lemma 8.** *There is a constant $c$ such that, for any $i \geq 1$ with $\lambda_i > 0$ and any $0 \leq k \leq n/c$, with probability at least $1 - 5e^{-n/c}$,*

$$\frac{\lambda_i^2 z_i^\top A_{-i}^{-2} z_i}{(1 + \lambda_i z_i^\top A_{-i}^{-1} z_i)^2} \geq \frac{1}{cn}\left(1 + \frac{\sum_{j>k} \lambda_j + n\lambda_{k+1}}{n\lambda_i}\right)^{-2}.$$

**Lemma 9.** *Suppose that $n \leq \infty$, $\{\eta_i\}_{i=1}^n$ is a sequence of non-negative random variables, and that $\{t_i\}_{i=1}^n$ is a sequence of nonnegative real numbers (at least one of which is strictly positive) such that, for some $\delta \in (0,1)$ and any $i \leq n$, $\Pr(\eta_i > t_i) \geq 1 - \delta$. Then,*

$$\Pr\left(\sum_{i=1}^n \eta_i \geq \frac{1}{2}\sum_{i=1}^n t_i\right) \geq 1 - 2\delta.$$

**Lemma 10.** *There are constants $c$ such that, for any $0 \leq k \leq n/c$ and any $b > 1$ with probability at least $1 - 10e^{-n/c}$,*

*1) if $r_k(\Sigma) < bn$, then $\mathrm{tr}(C) \geq \frac{k+1}{cb^2 n}$; and*

*2) if $r_k(\Sigma) \geq bn$, then*

$$\mathrm{tr}(C) \geq \frac{1}{cb^2}\min_{l \leq k}\left(\frac{l}{n} + \frac{b^2 n \sum_{i>l}\lambda_i^2}{(\lambda_{k+1} r_k(\Sigma))^2}\right).$$

*In particular, if all choices of $k \leq n/c$ give $r_k(\Sigma) < bn$, then $r_{n/c}(\Sigma) < bn$ implies that, with probability at least $1 - 10e^{-n/c}$, $\mathrm{tr}(C) = \Omega_{\sigma_r}(1)$.*

# 5.Proof

**Lemma 11.** *For any* $b \geq 1$ *and* $k^* := \min\{k : r_k(\Sigma) \geq bn\}$, *if* $k^* < \infty$, *we have*

$$
\min_{l \leq k^*} \left( \frac{l}{bn} + \frac{bn \sum_{i>l} \lambda_i^2}{(\lambda_{k^*+1} r_{k^*}(\Sigma))^2} \right)
$$

$$
= \frac{k^*}{bn} + \frac{bn \sum_{i>k^*} \lambda_i^2}{(\lambda_{k^*+1} r_{k^*}(\Sigma))^2} = \frac{k^*}{bn} + \frac{bn}{R_{k^*}(\Sigma)}.
$$

# 5.Proof

## • 5.4 Prove

**Lemma 10.** *There are constants $c$ such that, for any $0 \leq k \leq n/c$ and any $b > 1$ with probability at least $1 - 10e^{-n/c}$,*

1) *if $r_k(\Sigma) < bn$, then $\operatorname{tr}(C) \geq \frac{k+1}{cb^2n}$; and*
2) *if $r_k(\Sigma) \geq bn$, then*

$$\operatorname{tr}(C) \geq \frac{1}{cb^2} \min_{l \leq k} \left( \frac{l}{n} + \frac{b^2 n \sum_{i>l} \lambda_i^2}{(\lambda_{k+1} r_k(\Sigma))^2} \right).$$

**Upper Bound on the Trace Term. Lemma 7** *There are constants $b, c \geq 1$ such that, if $0 \leq k \leq n/c$, $r_k(\Sigma) \geq bn$, and $l \leq k$, then with probability at least $1 - 7e^{-n/c}$,*

$$\operatorname{tr}(C) \leq c \left( \frac{l}{n} + n \frac{\sum_{i>l} \lambda_i^2}{(\sum_{i>k} \lambda_i)^2} \right).$$

$$R(\hat{\theta}) \leq c \left( \|\theta^*\|^2 \, \|\Sigma\| \max \left\{ \sqrt{\frac{r_0(\Sigma)}{n}}, \frac{r_0(\Sigma)}{n}, \sqrt{\frac{\log(1/\delta)}{n}} \right\} \right.$$
$$\left. + c \log(1/\delta) \sigma_y^2 \left( \frac{k^*}{n} + \frac{n}{R_{k^*}(\Sigma)} \right) \right)$$

$$R(\hat{\theta}) \leq 2\theta^{*\top} B \theta^* + c\sigma^2 \log(1/\delta) \operatorname{tr}(C)$$

$$\mathbb{E} R(\hat{\theta}) \geq \frac{\sigma^2}{c} \left( \frac{k^*}{n} + \frac{n}{R_{k^*}(\Sigma)} \right).$$

$$\mathbb{E}_\varepsilon R(\hat{\theta}) \geq \theta^{*\top} B \theta^* + \sigma^2 \operatorname{tr}(C)$$

**Lemma 11.** *For any $b \geq 1$ and $k^* := \min\{k : r_k(\Sigma) \geq bn\}$, if $k^* < \infty$, we have*

$$\min_{l \leq k^*} \left( \frac{l}{bn} + \frac{bn \sum_{i>l} \lambda_i^2}{(\lambda_{k^*+1} r_{k^*}(\Sigma))^2} \right)$$

$$= \frac{k^*}{bn} + \frac{bn \sum_{i>k^*} \lambda_i^2}{(\lambda_{k^*+1} r_{k^*}(\Sigma))^2} = \frac{k^*}{bn} + \frac{bn}{R_{k^*}(\Sigma)}.$$

# 5.Proof

**Lemma 35.** *There is a constant $c$, that depends only on $\sigma_x$, such that for any $1 < t < n$, with probability at least $1 - e^{-t}$,*

$$\theta^{*\top} B\theta^* \leq c\|\theta^*\|^2\|\Sigma\| \max\left\{\sqrt{\frac{r_0(\Sigma)}{n}}, \frac{r_0(\Sigma)}{n}, \sqrt{\frac{t}{n}}\right\}.$$

# 6.Research Prospects and Future

**6.1 Conclusion:**

 1)We give finite sample excessrisk bounds that reveal the covariance structure that ensuresthat the minimum norm interpolating prediction rule has nearoptimal prediction accuracy.

 2)Overparameterization (that is, the existence of many **low-variance and hence**, unimportant directions in parameter space) is essential for benign overfitting and that data thatlie in a large but finite-dimensional space exhibit the benignoverfitting phenomenon with a much wider range of covariance properties than data that lie in an infinite-dimensionalspace.

# 6.Research Prospects and Future

**6.2 Future**

1）条件期望E($y|x$)不是$x$的一个线性函数

2）放宽"协变量作为独立随机变量向量的线性函数分布"这一条件

3）扩展损失函数

4）其他非线性参数化的函数类