

Statistics Learning Theory: Linear Regression

2019.03.02

Linear Regression

1. Domain set $\mathcal{X} \in R^d$
2. Label set \mathcal{Y} is the set of real number
3. Hypothesis class

$$\mathcal{H}_{reg} = L_d = \{x \mapsto \langle w, x \rangle + b : w \in R^d, b \in R\}$$

4. loss function

$$l(h, (x, y)) = (h(x) - y)^2$$

5. Empirical Risk function

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

Linear Regression: Implementation of ERM rule

Least square is the algorithm that solves the ERM problem for hypothesis class of linear regression predictors with respect to the squared loss.

The ERM problem can be written as

$$\arg \min_w L_S(h_w) = \arg \min_w \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

Linear Regression: Closed form solution for ERM rule

To solve the problem we calculate the gradient of the objective function

$$\frac{2}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i) x_i = 0$$

We can rewrite the problem as the problem $Aw = b$ where

$$A = \left(\sum_{i=1}^m x_i x_i^T \right) \quad \text{and} \quad b = \sum_{i=1}^m y_i x_i$$

If A is invertible then the solution to the ERM problem is

$$w = A^{-1}b$$

Linear Regression: ERM rule

If A is not invertible, we still can find the solution to the system $Aw = b$ since b is in the range of A .

To be specific, since A is symmetric we can decompose A as

$$A = VDV^T$$

where D is a diagonal matrix and V is an orthonormal matrix. After normalizing D we can obtain

$$A^+ = VD^+V^T \quad \text{and} \quad \hat{w} = A^+b$$

Let v_i denote the i 'th column of V . Then, we have

$$A\hat{w} = AA^+b = VDV^TVD^+V^Tb = VDD^+V^Tb = \sum_{i:D_{i,i} \neq 0} v_i v_i^T b$$

Consistency Theorem

Let

$Q_n(\theta)$ = some objective function

$$\hat{\theta} = \arg \max_{\theta \in \Theta} Q_n(\theta)$$

Examples:

1. NLLS : $Q_n(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i - m(x_i, \theta))^2$
2. ML : $Q_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(z_i, \theta)$
3. GMM : $Q_n(\theta) = -(\frac{1}{n} \sum_{i=1}^n g(z_i, \theta))' W (\frac{1}{n} \sum_{i=1}^n g(z_i, \theta))$

Consistency Theorem

Theorem (General Consistency Theorem)

Suppose

1. Θ is compact
2. $\sup_{\theta \in \Theta} |Q_n(\theta) - Q_*(\theta)| \rightarrow^P 0$ for some $Q_* : \Theta \rightarrow R$
3. Q_* is continuous in $\theta \in \Theta$
4. Q_* is uniquely maximized at θ_0

Then

$$\hat{\theta} \rightarrow^P \theta_0$$

Consistency Theorem Proof

Pick $\varepsilon > 0$, since $\hat{\theta}$ maximizes $Q_n(\theta)$

$$Q_n(\hat{\theta}) > Q_n(\theta_0) - \frac{\varepsilon}{3}$$

By condition 2, for any $\theta \in \Theta$

$$|Q_n(\theta) - Q_*(\theta)| < \frac{\varepsilon}{3}$$

with probability approaching one

Consistency Theorem Proof

Thus with probability approaching one

$$Q_n(\hat{\theta}) - Q_*(\hat{\theta}) < \frac{\varepsilon}{3}$$

$$Q_*(\theta_0) - Q_n(\theta_0) < \frac{\varepsilon}{3}$$

Combining these inequalities

$$Q_*(\hat{\theta}) + \frac{\varepsilon}{3} > Q_n(\hat{\theta})$$

$$> Q_n(\theta_0) - \frac{\varepsilon}{3}$$

$$> Q_*(\theta_0) - \frac{2\varepsilon}{3}$$

Therefore,

$$Q_*(\hat{\theta}) > Q_*(\theta_0) - \varepsilon$$

Consistency Theorem Proof

Since we want to prove that $\hat{\theta} \rightarrow^P \theta_0$, we want to show that

$$Pr\{\hat{\theta} \in \mathcal{N}\} \rightarrow 1$$

for any open set $\mathcal{N} \subset \Theta$ containing θ_0

Now pick any open set \mathcal{N} containing θ_0 , since \mathcal{N} is open, \mathcal{N}^c is closed. By condition 1, Θ is compact, $\Theta \cap \mathcal{N}^c$ is also compact. Since Q_* is continuous, Weierstrass theorem guarantees there exists $\theta_* \in \Theta \cap \mathcal{N}^c$ such that

$$\sup_{\Theta \cap \mathcal{N}^c} Q_*(\theta) = Q_*(\theta_*)$$

Consistency Theorem Proof

Since Q_* is uniquely maximized at θ_0 , we have

$$Q_*(\theta_0) > Q_*(\theta_*)$$

and set

$$\varepsilon' = Q_*(\theta_0) - Q_*(\theta_*) > 0$$

Using the previous inequality and set $\varepsilon = \varepsilon'$

$$\begin{aligned} Q_*(\hat{\theta}) &> Q_*(\theta_0) - \varepsilon' \\ &= Q_*(\theta_*) \\ &= \sup_{\Theta \cap \mathcal{N}^c} Q_*(\theta) \end{aligned}$$

This means that

$$\hat{\theta} \in \mathcal{N}$$

with probability approaching one

Consistency Theorem Remark

1. For each application, most efforts are devoted to check Conditions 2 and 4
2. Condition 4 is called identification condition. If $Q_*(\theta)$ is maximized at multiple points, we cannot tell where $\hat{\theta}$ converges in general
3. Condition 2 says that objective function $Q_n(\theta)$ uniformly converges in probability to the limit objective function $Q_*(\theta)$
4. For condition 2, we typically need some kind of uniform law of large numbers

Uniform law of large number

Theorem

1. Θ is compact
2. $g(z, \theta)$ is almost surely continuous at each $\theta \in \Theta$
3. There is $d(z)$ such that $|g(z, \theta)| \leq d(z)$ for all $\theta \in \Theta$ and almost every z and $E[d(z)] < \infty$

Then

$$\sup_{\theta \in \Theta} |\bar{g}(\theta) - E[g(z, \theta)]| \xrightarrow{P} 0$$

Consistency of NLLSE

Theorem

Suppose

1. $\{(y_i, x_i')\}_{i=1}^n$ is iid and $E[y|x] = m(x, \theta)$ almost surely only if $\theta = \theta_0$
2. Θ is compact
3. $m(x, \theta)$ is almost surely continuous at each $\theta \in \Theta$
4. $E[y^2] < \infty$ and $E[\sup_{\theta \in \Theta} |m(x, \theta)|^2] < \infty$

Then

$$\hat{\theta} \rightarrow^p \theta_0$$

Consistency of NLLSE

It is sufficient to check condition 1 – 4 in the general theorem.

Condition 1 is guaranteed by our second condition.

Condition 2: $Q_*(\theta) = -E[\{y - m(x, \theta)\}^2]$, we want to show that

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \{y_i - m(x_i, \theta)\}^2 - E[\{y - m(x, \theta)\}^2] \right| \rightarrow^p 0$$

The above is holded by applying the ULLN. Since ULLN also guarantees continuity of $Q_*(\theta)$. Thus condition 3 is also satisfied.

Consistency of NLLSE

It remains to check condition 4, identification of θ_0 . i.e.

$$Q_*(\theta) = -E[\{y - m(x, \theta)\}^2]$$

is uniquely maximized at θ_0 . Since $m(x) = E[y|x]$ solves

$$\min_g E[\{y - g(x)\}^2]$$

General Asymptotic Normality Theorem: Basic Idea

Now consider asymptotic distribution of extremum estimator

$Q_n(\theta)$ = some objective function

$$\hat{\theta} = \arg \max_{\theta \in \Theta} Q_n(\theta)$$

Assume consistency $\hat{\theta} \rightarrow^P \theta_0$

We want to derive asymptotic normal distribution in the form of

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, V)$$

The result can be used to construct confidence interval or to conduct hypothesis testing

General Asymptotic Normality Theorem: Basic Idea

Suppose $Q_n(\theta)$ is continuously twice differentiable. Look at FOC for $\hat{\theta}$

$$\frac{\partial Q_n(\hat{\theta})}{\partial \theta} = 0$$

Using the Taylor expansion, we can get

$$0 = \frac{\partial Q_n(\theta_0)}{\partial \theta} + \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} (\hat{\theta} - \theta_0)$$

where $\tilde{\theta}$ is a point on the line joining $\hat{\theta}$ and θ_0 . Solving for $(\hat{\theta} - \theta_0)$ (if inverse exists), we get

$$\sqrt{n}(\hat{\theta} - \theta_0) = -\left(\frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'}\right)^{-1} \left(\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta}\right)$$

General Asymptotic Normality Theorem: Basic Idea

If

$$\frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} \rightarrow^p H$$
$$\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \rightarrow^d N(0, \Sigma)$$

Then we obtain

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, H^{-1} \Sigma H^{-1})$$

Typically the first condition is verified by ULLN and the second condition is verified by CLT.

General Asymptotic Normality Theorem

Theorem

Suppose

1. $\hat{\theta} \rightarrow^P \theta_0$ and $\theta_0 \in \text{int}\Theta$
2. $Q_n(\theta)$ is twice continuously differentiable in a neighborhood \mathcal{N} of θ_0
3. $\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \rightarrow^d N(0, \Sigma)$
4. There exists $H(\theta)$ that is continuous at θ_0 ,
 $\sup_{\theta \in \mathcal{N}} \left| \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} - H(\theta) \right| \rightarrow^P 0$, and $H = H(\theta_0)$ is nonsingular

Then

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, H^{-1}\Sigma H^{-1})$$

General Asymptotic Normality Theorem: Proof

Pick any convex and open set $\mathcal{N}' \subset \mathcal{N}$ containing θ_0 , and define the indicator $\hat{I} = I\{\hat{\theta} \in \mathcal{N}'\}$

Note that $\hat{I} \xrightarrow{P} 1$

By condition 2 and Taylor expansion we get

$$\begin{aligned} 0 &= \hat{I} \frac{\partial Q_n(\hat{\theta})}{\partial \theta} \\ &= \hat{I} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + \hat{I} \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} (\hat{\theta} - \theta_0) \end{aligned}$$

General Asymptotic Normality Theorem: Proof

Since $\tilde{\theta} \rightarrow^P \theta_0$, we have

$$\begin{aligned} & \left| \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} - H \right| \\ & \leq \left| \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} - H(\tilde{\theta}) \right| + |H(\tilde{\theta}) - H| \\ & \leq \sup_{\theta \in \mathcal{N}} \left| \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} - H(\theta) \right| + |H(\tilde{\theta}) - H| \\ & \rightarrow^P 0 \end{aligned}$$

By condition 4 and continuous mapping theorem

Since H is nonsingular,

$$\bar{I} := I\left\{\hat{\theta} \in \mathcal{N}', \frac{\partial^2 Q_n(\tilde{\theta})}{\partial \theta \partial \theta'} \text{ is nonsingular}\right\} \rightarrow^P 1$$

General Asymptotic Normality Theorem: Proof

Combining these results,

$$0 = \bar{I}\sqrt{n}\frac{\partial Q_n(\theta_0)}{\partial\theta} + \bar{I}\frac{\partial^2 Q_n(\tilde{\theta})}{\partial\theta\partial\theta'}(\hat{\theta} - \theta_0)$$

Solving for $(\hat{\theta} - \theta_0)$ we obtain

$$\begin{aligned} & \sqrt{n}(\hat{\theta} - \theta_0) \\ &= -\bar{I}\left(\frac{\partial^2 Q_n(\tilde{\theta})}{\partial\theta\partial\theta'}\right)^{-1}\left(\sqrt{n}\frac{\partial Q_n(\theta_0)}{\partial\theta}\right) + \{1 - \bar{I}\}\sqrt{n}(\hat{\theta} - \theta_0) \end{aligned}$$

By condition 3 and $\bar{I} \rightarrow^p 1$, the first term converges to $N(0, H^{-1}\Sigma H^{-1})$.

General Asymptotic Normality Theorem: Proof

It remains to show second term satisfies

$$\{1 - \bar{I}\}\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^P 0$$

To see this, let $Z_n = \sqrt{n}(\bar{\theta} - \theta_0)$ and pick any $\varepsilon > 0$. Then

$$\begin{aligned} & Pr\{|(1 - \bar{I}Z_n| > \varepsilon\} \\ &= Pr\{|(1 - \bar{I}Z_n| > \varepsilon, \bar{I} = 1\} + Pr\{|(1 - \bar{I})Z_n| > \varepsilon, \bar{I} = 0\} \\ &\leq Pr\{0 > \varepsilon\} + Pr\{\bar{I} = 0\} \\ &\rightarrow 0 \end{aligned}$$

General Asymptotic Normality Theorem: Application to NLLSE

Let

$$Q_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \{y_i - m(x_i, \theta)\}^2$$

$$\hat{\theta} = \arg \max_{\theta \in \Theta} Q_n(\theta)$$

General Asymptotic Normality Theorem: Application to NLLSE

Theorem

Suppose

1. $\theta_0 \in \text{int}\Theta$
2. $m(x, \theta)$ is twice continuously differentiable in a neighborhood \mathcal{N} of θ_0 with probability one
3. $E[y^4] < \infty$ $E\left|\frac{\partial m(x, \theta_0)}{\partial \theta}\right|^4$ and

$$E\left[\sup_{\theta \in \mathcal{N}} \left|\frac{\partial^2 m(x, \theta)}{\partial \theta \partial \theta'}\right|^2\right] < \infty$$

4. $H = E\left[\frac{\partial m(x, \theta_0)}{\partial \theta} \frac{\partial m(x, \theta_0)}{\partial \theta'}\right]$ is nonsingular

Then

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, H^{-1}\Sigma H^{-1})$$

where $\Sigma = E\left[e^2 \frac{\partial m(x, \theta_0)}{\partial \theta} \frac{\partial m(x, \theta_0)}{\partial \theta'}\right]$ and $e = y - m(x, \theta_0)$

General Asymptotic Normality Theorem: Application to NLLSE Proof

It is enough to check conditions for General asymptotic normality theorem.

For condition 1, consistency is already verified and $\theta_0 \in \text{int}\Theta$

Condition 2 is satisfied by (ii)

In this case, condition 3 is verified as

$$\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial m(x_i, \theta_0)}{\partial \theta} e_i \rightarrow^d N(0, \Sigma)$$

where $e_i = y_i - m(x, \theta_0)$ and convergence follows from CLT with (iii) (Note that $E[y^4] < \infty \rightarrow E[e^4] < \infty$)