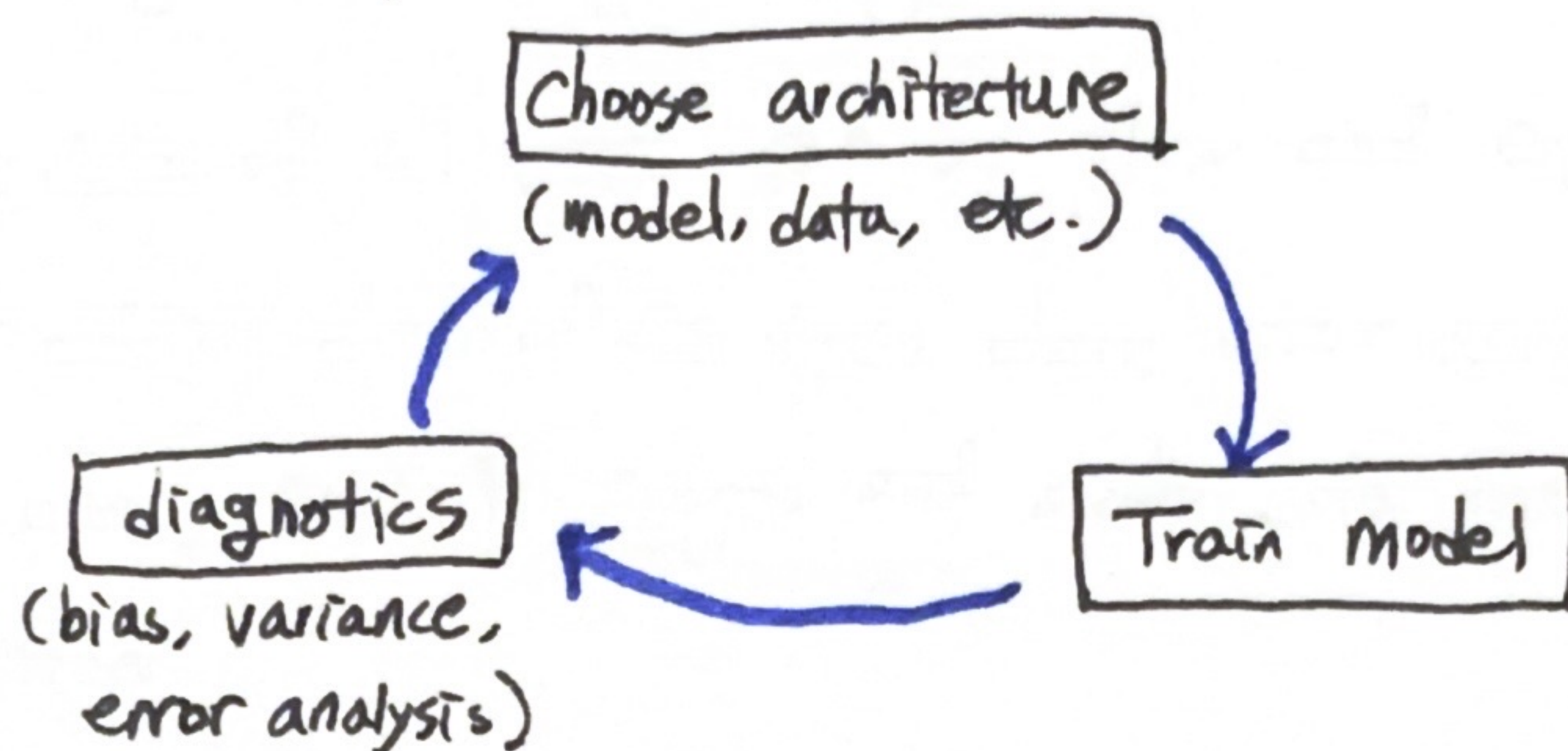


<Machine Learning development process - Iterative loop of ML development>

* Iterative loop of ML development



ex) Spam mail classifier

Supervised learning : \vec{x} = features of email list the top 10,000 words to compute $x_1, x_2, \dots, x_{10,000}$
 y = spam(1) or not spam(0) then \vec{x} vector is : 추출한 10,000개 단어에 mail에 포함되었는지 여부로 feature vector

↓ after initializing a learning algorithm
how to try to reduce spam classifier algorithm's error?

- ① Collect more data
- ② Develop sophisticated features based on email routing (from email header)
- ③ Define sophisticated features from email body (e.g. should "discounting" and "discount" be treated as same word)
- ④ Design algorithms to detect misspellings (e.g. watches, medicine, mortgage)

<Error Analysis>

- cross validation set 중 misclassified된 데이터 example들을 수동으로 분석하여 어떤 type의 데이터가 어떤 점량의 error를 발생시키는 지 관찰

ex) $M_{cv} = 500$ examples in CV set.
Algorithm misclassified 100 of them

⇒ Manually examine 100 examples and categorize them based on common traits

ex) - Pharma : 21

- Deliberate misspelling : 3
- Unusual email routing : 11
- Steal password (phishing) : 18
- Spam message in embedded image : 5

→ 통계적으로 misclassified된 example의 category에 대한 상을 붙이는 방법으로 next step을 결정

- Pharma spam email collect more data about pharma
come up with new features
(e.g. specific name of drug, pharma product)
- Phishing email collect more data about phishing email
come up with new features
(e.g. specific url in email...)