# Applied ML is a highly iterative process

# layers

# hidden units

learning rates

activation functions

...

Idea

Experiment

Code

NLP, Vision, Speech, Structural data

Ads    Search    Security    logistic ....

Andrew Ng

# Train/dev/test sets

Data

train set

dev    test

- Hold-out
  cross validation
- Development set
  "dev"

test

Previous era :    70/30 %

previous era

100 - 1000 - 10000

60/20/20 %

train    test

Big data :    1,000,000    10,000    10,000

Big data era

98 / 1 / 1 %

99.5 / .25 / .25
         / .4 / .1 %

전통적인 방법론에서 사용하는것처럼 10~20%
씩이나 dev/test set에 사용할 필요는 없음

Andrew Ng

# Mismatched train/test distribution

Cats

Training set:
Cat pictures from
webpages

Dev/test sets:
Cat pictures from
users using your app

Make sure dev and test come from Same distribution.

"test"

train /dev

train /test

→ Tron /dev

Not having a test set might be okay. (Only dev set.)

Andrew Ng