# Developing and Evaluating an Anomaly Detection System

ex. Aircraft engines motivating ex.

- 10000 good (normal) engines
- 20 flawed (anomalous) engines

→ Training set: 6000 good engines $(y=0)$

CV: 2000 good engines $(y=0)$  10 anomalous $(y=1)$

Test: 2000 good engines $(y=0)$  10 anomalous $(y=1)$

## Algorithm evaluation

Fit model $p(x)$ on training set $\{x^{(1)}, \ldots, x^{(m)}\}$

On a cv/test example $x$, predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$
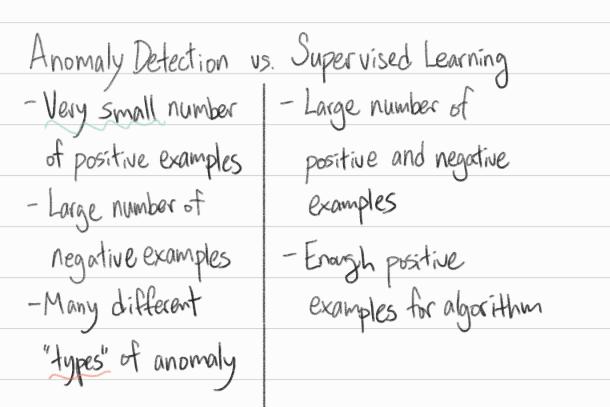
skew data인 경우, accuracy는 부적절

so, TP·FP·FN·TN

Precision / Recall

$F_1$-score

(can also use cross validation set to choose parameter $\varepsilon$)

# Anomaly Detection vs. Supervised Learning

| | |
|---|---|
| - Very small number of positive examples | - Large number of positive and negative examples |
| - Large number of negative examples | - Enough positive examples for algorithm |
| - Many different "types" of anomaly | |

# Choosing What Features to Use

for Non-gaussian features



$$x \rightarrow \log(x + C)$$
$$x \rightarrow \sqrt{x}$$

Gaussian 형태로 변환

Error analysis for anomaly detection

Want $p(x)$ large for $0$

$p(x)$ small for $1$

Most common prob.: $p(x)$ is comparable for $0$ and $1$