# Classification metrics optimization: Logloss and accuracy

# Classification metrics optimization

- Logloss

- Accuracy

- AUC

- (Quadratic weighted) Kappa

# Logloss

$$\text{LogLoss} = -\frac{1}{N} \sum_{i=1}^{N} y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$

How do you optimize it?

Just run the right model!
(or calibrate others)

# Logloss

- **Tree-based**

  `XGBoost, LightGBM`

  ~~`sklearn.RandomForestClassifier`~~

- **Linear models**

  `sklearn.<>Regression`

  `sklearn.SGDRegressor`

  `Vowpal Wabbit`

- **Neural nets**

  `PyTorch, Keras, TF, etc.`

*Synonyms: Logistic loss*                        Read the docs!

# Logloss

**분류기에서 계산한 확률과 실제 데이터의 확률이 다를 수 있음, 왜?**
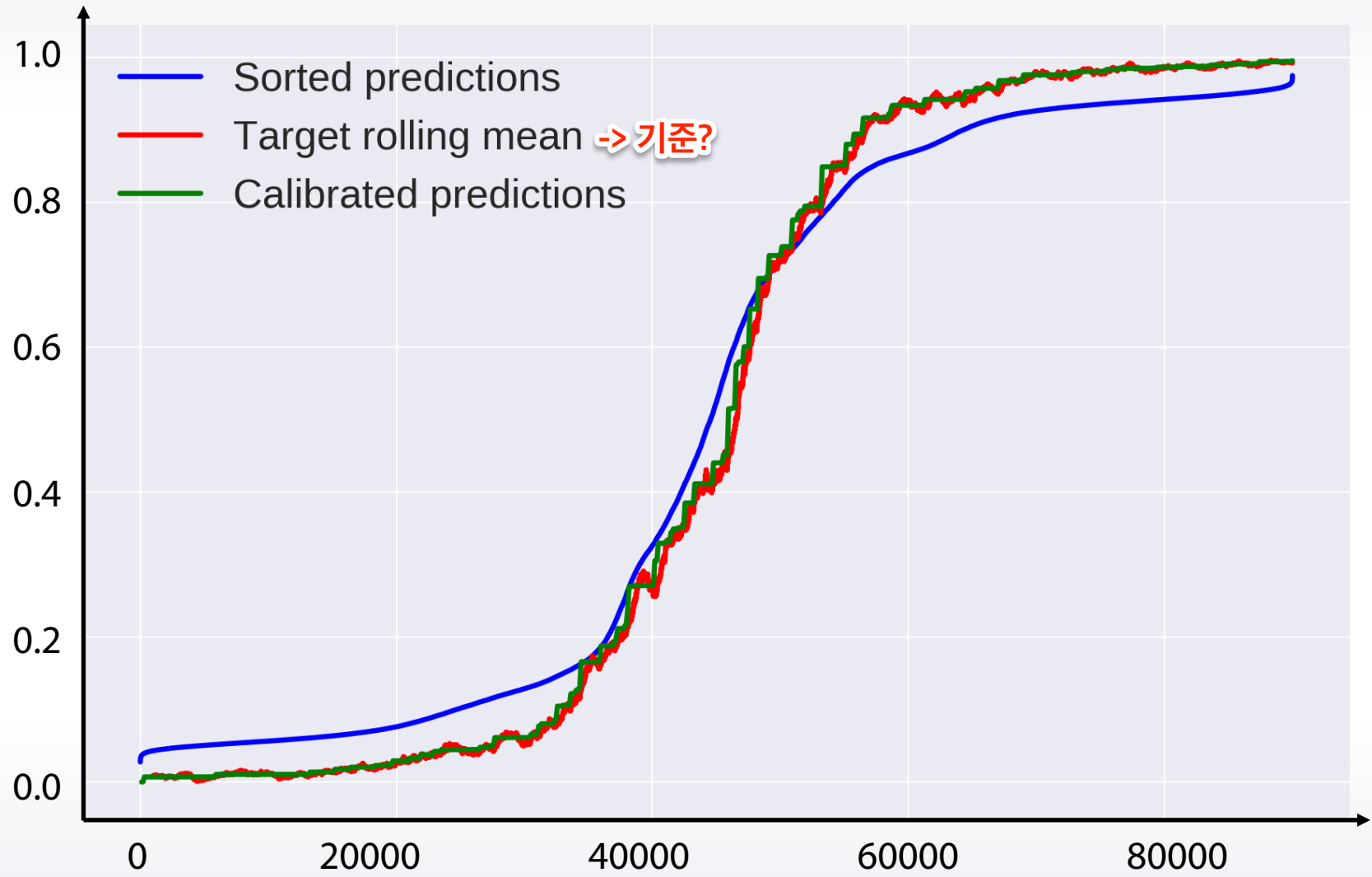
Correct probabilities:

- Take all objects with score e.g. ~ 0.8
    - 80% of them of class 1
    - 20% of them class 0


Incorrect probabilities:

- Take all objects with score e.g. ~ 0.8
    - 50% of them of class 1
    - 50% of them of class 0

# Probability calibration 보정하면 좋다



Sorted predictions
Target rolling mean -> 기준?
Calibrated predictions

# Probability calibration

logistic 외에 svm, decision tree에서도 예측 확률을 돌려주는 방법들이 있음
확률로 바꾼 것을 여러번 cross validation하여 보정함(sklearn.calibration의 CalibratedClassifierCV)

- Platt scaling

  – Just fit Logistic Regression to your predictions
    (like in stacking) classifier(ex: SVM)의 output을 확률로 변환
    로지스틱 공식이랑 비슷, 시그모이드 형태에 적합

- Isotonic regression

  – Just fit Isotonic Regression to your predictions
    (like in stacking) monotinic regression, non-parametric model

- Stacking

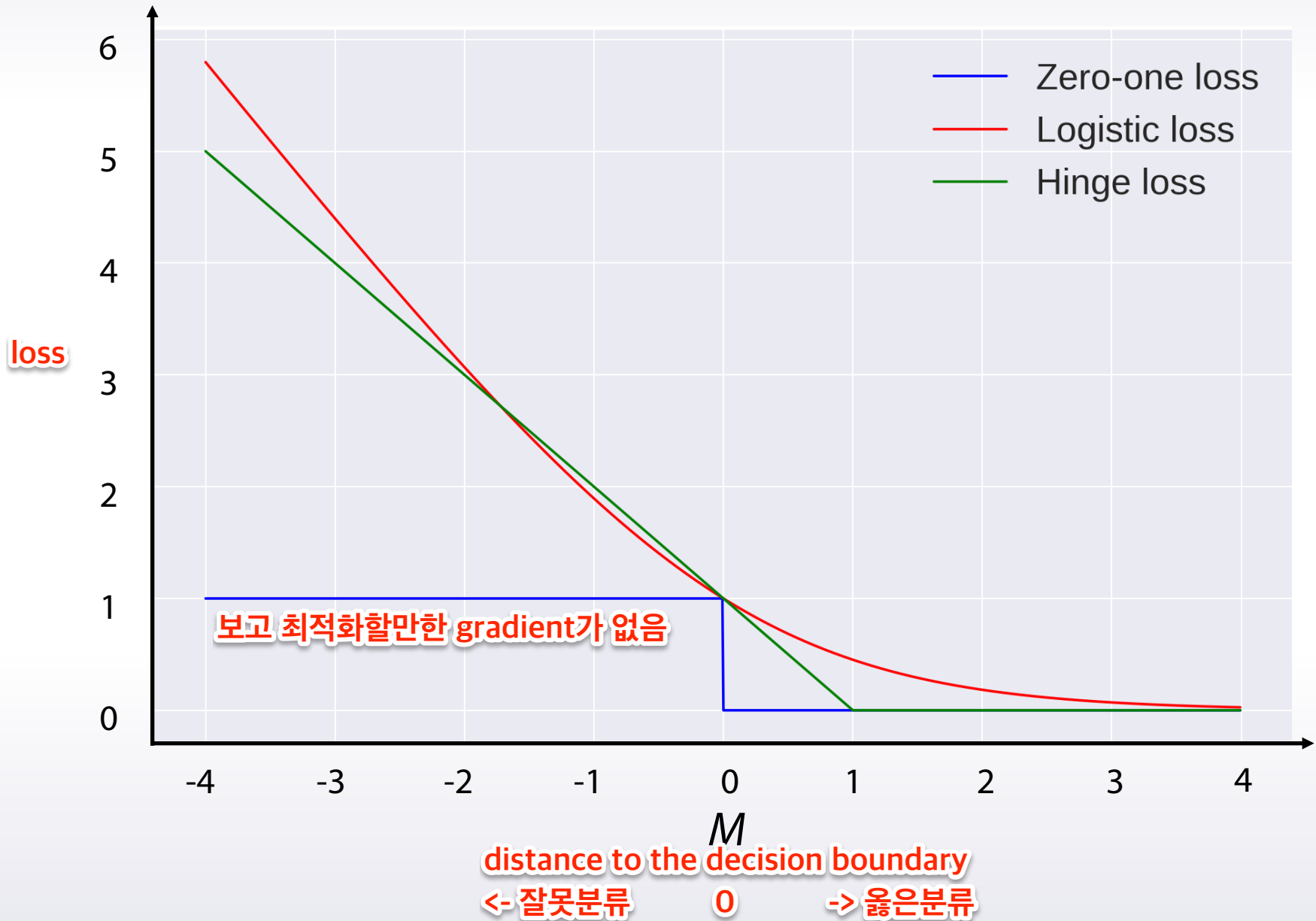  – Just fit XGBoost or neural net to your predictions

# Accuracy

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^{N} [\hat{y}_i = y_i]$$

How do you optimize it?
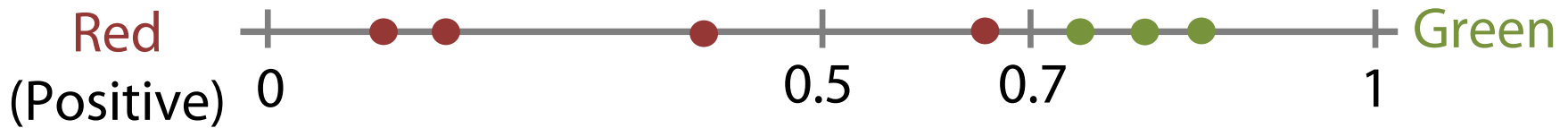
Fit any metric and tune treshold!

# Accuracy

Accuracy를 최적화하기 어려운 이유

- Zero-one loss
- Logistic loss
- Hinge loss

loss

보고 최적화할만한 gradient가 없음

$M$

distance to the decision boundary

<- 잘못분류　　　0　　　-> 옳은분류

# Accuracy

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^{N} [[f(x) > b] = y_i]$$

Red
(Positive)



Green

0    0.5    0.7    1

$$b = 0.5 \quad \Rightarrow \quad \text{Accuracy} = \frac{6}{7}$$

$$b = 0.7 \quad \Rightarrow \quad \text{Accuracy} = 1$$

# Conclusion

- Logloss

- Accuracy

- AUC

- (Quadratic weighted) Kappa