

Question 1

Why can not we directly use outputs of linear transformation for classification?

1. Scale of original outputs are **not lying in range of 0 to 1**, which does not meet the limitation of possibility

2. Hard-code classificaiton like  $\{0, 1\}$  is **not suitable for derivative**, so is hard for optimization

3. Original outputs may **include negative values**

$$f_{\theta}(x) = \sigma(\theta^T x)$$

Logistic (sigmoid)

Softmax

when  $K = 2$  get

Transformation Algorithms

Advantage

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

1. symmetry at point (0, 0.5)  
2. sensitive to change of x  
3. continuous, monotonically increasing, derivable

Transformation Algorithms

Advantage

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

1. not negative output  
2. suitable for multi-classify  
3. sum up each equals to 1

Question 2

How can we choose loss functions? And How can we understand them?

Persepctive of statistic/possibility

Maximum likelihood estimation

$$L(\theta) = \prod_{i=1}^N f_{\theta}(x_i)^{y_i} (1 - f_{\theta}(x_i))^{1-y_i}$$

Production is hard to optimize and derive

Log-likelihood

$$l(\theta) = \log L(\theta) = \sum_{i=1}^N [y_i \log f_{\theta}(x_i) + (1 - y_i) \log(1 - f_{\theta}(x_i))]$$

Summary is easy to optimize and it is monotonically increasing

log

~ the same to minus ~

Persepctive of entropy

Cross entropy

$$\begin{aligned} H(p, q) &= - \sum_{i=1}^n p(X_i) \log q(X_i) \\ &= -p(X=1) \log q(X=1) - p(X=0) \log q(X=0) \\ &= -y \log f_{\theta}(x) - (1-y) \log(1 - f_{\theta}(x)) \end{aligned}$$

To minimize relative entropy is to minimize  $H(p, q)$

relative entropy

$$\begin{aligned} D_{KL}(p \parallel q) &= E_{X \sim p(x)} \left[ \log \left( \frac{p(x)}{q(x)} \right) \right] \\ &= \sum_{i=1}^n (p(X_i) \log p(X_i) - p(X_i) \log q(X_i)) \\ &= -H(p) + H(p, q) \end{aligned}$$

KL divergence = average distance between p and q distribution with p possibility as weights

Information/Entropy

information

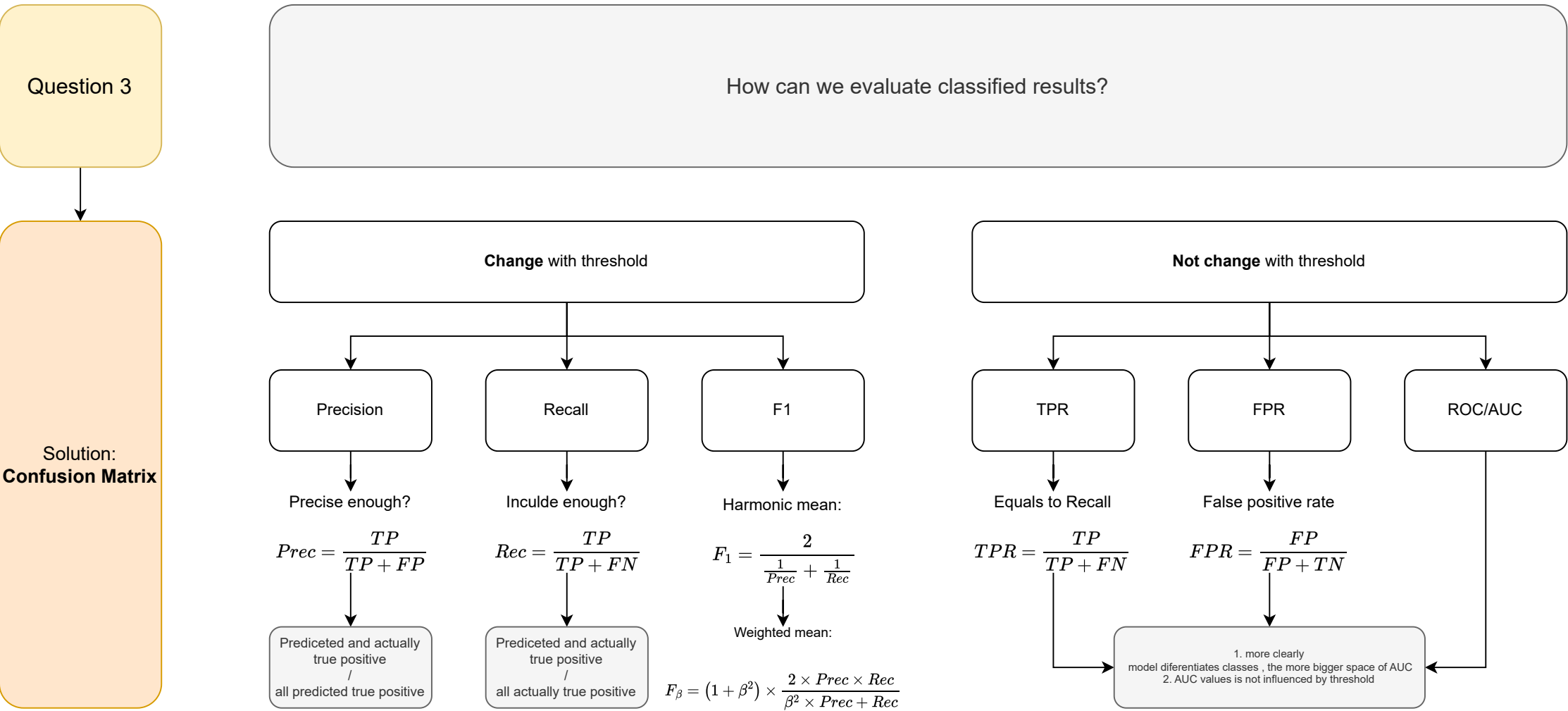
$$I(X_i) = -\log P(X_i)$$

Entropy

$$H(p) = E_{X \sim p(x)} [I(X)] = \sum_{i=1}^n P(X_i) I(X_i) = - \sum_{i=1}^n P(X_i) \log P(X_i)$$

Distance between p and q distribution

ic and Softmax transformation



GLM, Generalized Linear Model

