

# Evaluating predict

Monday, April 24, 2017 3:25 PM

Have test set  $\{x_i, y_i\}$ , want to evaluate

predict method  $\hat{y}_i = \begin{cases} 1, & s_i > \tau \\ 0, & s_i \leq \tau \end{cases}$

Confusion Matrix

	Predict 1	Predict 0
Actual 1	True Positive (TP)	False Negative (FN)
Actual 0	False Positive (FP)	True Negative (TN)

Predict

1. Calculate score  $s_i$  for each  $i$ .

2. Order scores :  $s_{a_1} \geq s_{a_2} \geq s_{a_3} \geq \dots \geq s_{a_n}$

where  $a_1, \dots, a_n$  is permutation of  $\{1, \dots, n\}$ .

3. Predict  $\hat{y}_{a_1}, \dots, \hat{y}_{a_T} = 1$  and  $\hat{y}_{a_{T+1}}, \dots, \hat{y}_{a_n} = 0$ .

$$S : 3.1 \ 3 \ 2.8 \ 2.6 \ 2.3 \ 2.1 \ 1.1 \ 1.8 \ .6 \ .5 \ .2$$

$$\hat{y} \quad | \quad 1 \ \ 1 \ \ 1 \ \ 1 \ \ 1 \ \ 1 \ \ 0 \ \ 0 \ \ 0 \ \ 0 \ \ 0 \ \ 0 \quad T=6$$

$$y \quad | \quad 1 \ \ 0 \ \ 1 \ \ 1 \ \ 0 \ \ 1 \ \ 0 \ \ 1 \ \ 0 \ \ 0 \ \ 1 \ \ 0$$

$$\text{Conf.} \quad TP - 4 \quad | \quad FN = 7$$

Conf.  
Matrix

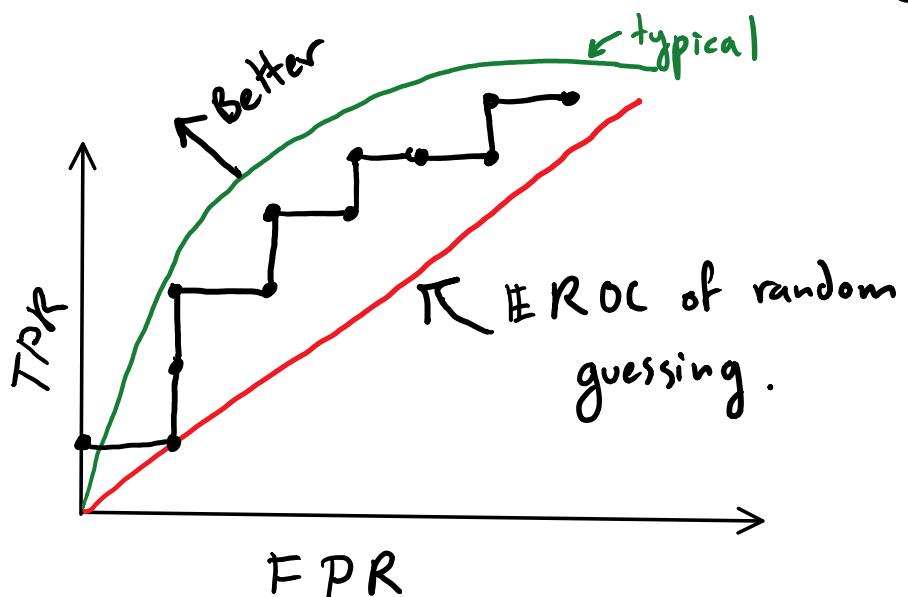
$TP = 4$	$FN = 2$
$FP = 2$	$TN = 4$

## Reciever - Operating Characteristic (ROC)

True positive rate (TPR) =  $\frac{TP}{TP + FN}$  } Actual Positive  
(aka recall)

False positive rate (FPR) =  $\frac{FP}{FP + TN}$  } Actual Neg.

T	1	2	3	4	5	6	7	8
TPR	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{3}{6}$	$\frac{3}{6}$	$\frac{4}{6}$	$\frac{4}{6}$	$\frac{5}{6}$
FPR	$\frac{0}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{2}{6}$	$\frac{3}{6}$	$\frac{3}{6}$



- Every element of confusion matrix is represented in ROC

## Precision - Recall Curve

In recommendation systems, many 0's - few 1's.

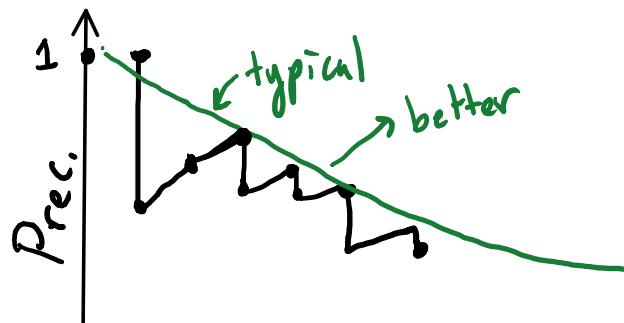
- Ex
- Link prediction in sparse graphs (friend recom.)
  - Search engine (most pages are not relevant)
  - Document retrieval

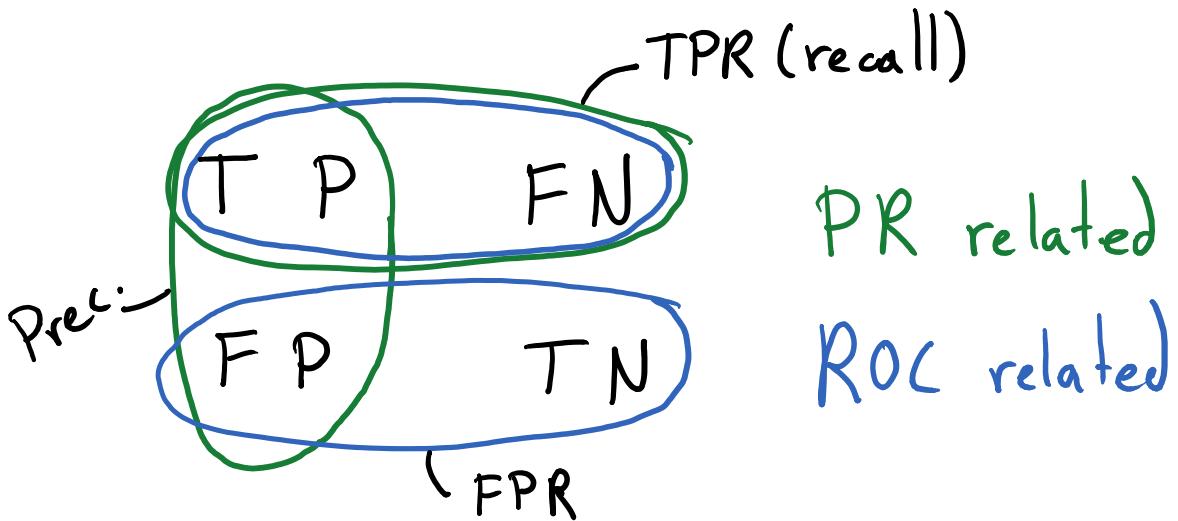
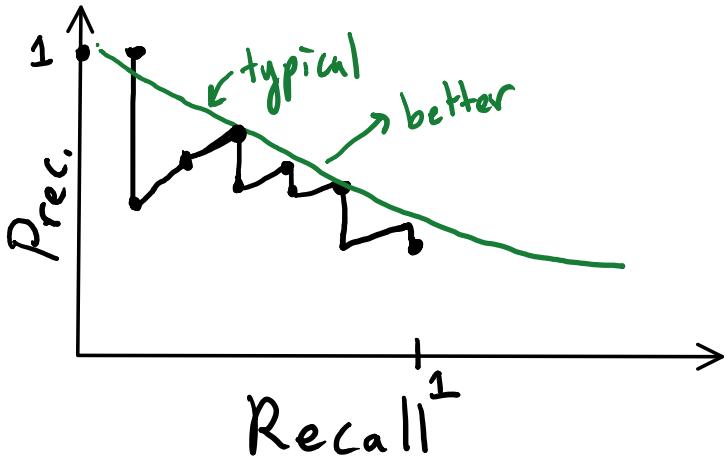
$$\text{Recall} = TPR = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

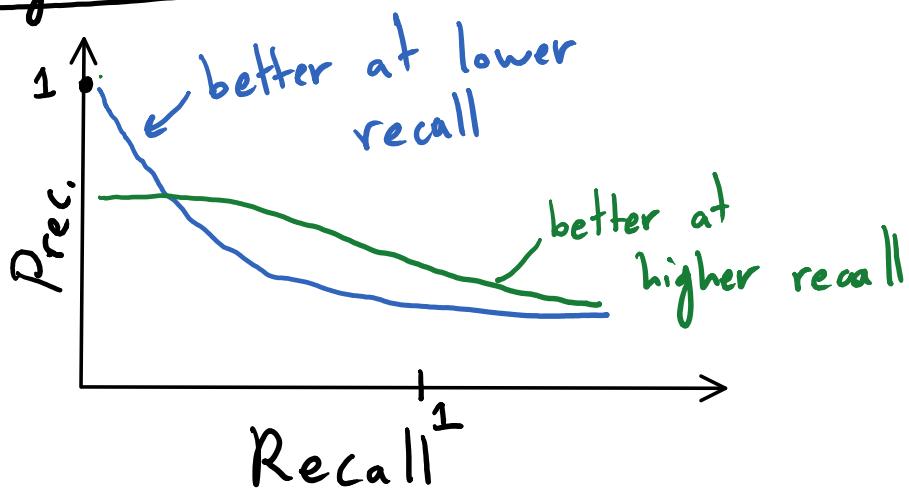
} Neither is sensitive to increasing TN

T	1	2	3	4	5	6	7	8
Recall	1/6	1/6	2/6	3/6	3/6	4/6	4/6	5/6
Prec.	1	1/2	2/3	3/4	3/5	4/6	4/7	5/8





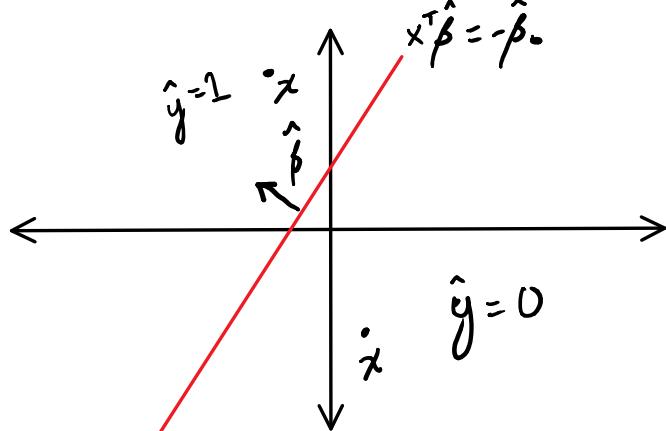
### Comparing PR Curves



# Margin Based Methods

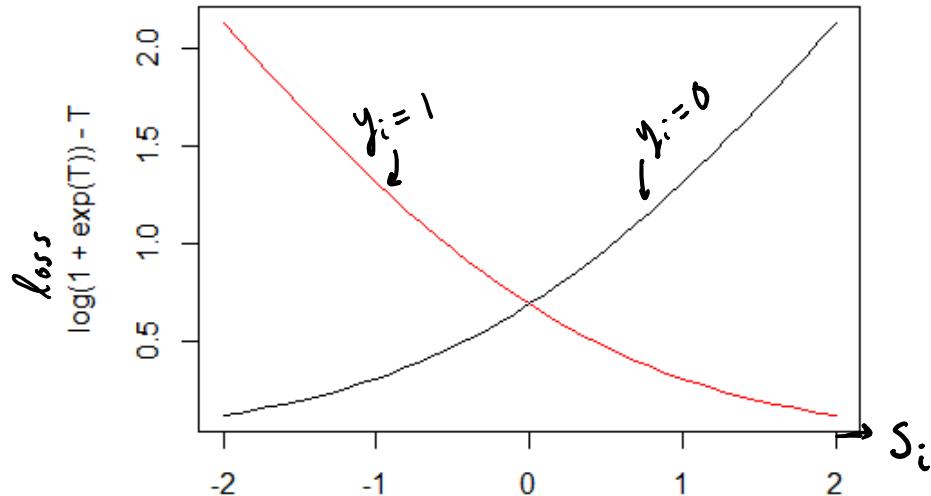
Tuesday, April 25, 2017 4:22 PM

$$\hat{y} = \begin{cases} 1, & \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p \geq 0 \\ 0, & \text{otherwise} \end{cases}$$



Empirical Risk Minimization

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \text{loss}(y_i, x_i, \beta)$$



$$s_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

$$\underline{0-1 \text{ loss}} (y_i = 1)$$

$$\text{loss}_{0-1} = \begin{cases} 1, & s_i < 0 \\ 0, & \text{otherwise} \end{cases}$$

logistic regression ( $y_i = 1$ )

$$\min \frac{1}{n} \sum_i \text{loss}_{0-1}(y_i, x_i, \beta)$$

is very hard to optimize!

Support vector machines

logistic regression ( $y_i=1$ )

$$\text{loss} = \log(1+e^{-s_i})$$

(for  $y_i=0$  switch)  
 $s_i \leftarrow -s_i$

Support vector machines

$$\text{loss} = \begin{cases} 1-s_i & , s_i < 1 \\ 0 & , s_i \geq 1 \end{cases}$$

$$\min_{\beta} \frac{1}{n} \left[ \sum_{i:y_i=1} \text{loss}(1, x_i, \beta) + \sum_{i:y_i=0} \text{loss}(0, x_i, \beta) \right]$$

$$\min_{\beta} \frac{1}{n} \left[ \sum_{i:y_i=1} \alpha \text{loss}_{0,1}(1, x_i, \beta) + \sum_{i:y_i=0} \text{loss}(0, x_i, \beta) \right]$$

$$\alpha = \frac{\sum_{i:y_i=1} 1}{n}$$

# Overfitting and the Kernel Trick

Friday, May 12, 2017 11:57 AM

$x_i$  - 2 dimensional feature vector

$\varphi_i(x_i)$  -  $\lambda^{\text{th}}$  feature in high-dimensional embedding

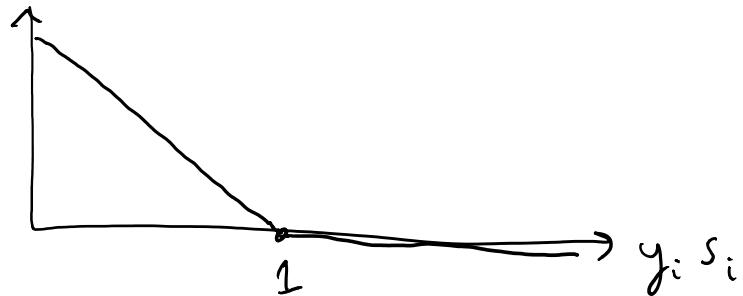
ex  $\varphi_1(x_i) = 1, \varphi_2(x_i) = x_{i1}, \varphi_3(x_i) = x_{i2},$   
 $\varphi_4(x_i) = x_{i1}^2, \varphi_5(x_i) = x_{i1} \cdot x_{i2}, \varphi_6(x_i) = x_{i2}^2$

$z_i = \varphi(x_i)$  - 6 dimensional engineered feature vector

Support Vector Machines

$$\min_{\beta \in \mathbb{R}^P} \frac{1}{n} \sum_{i=1}^n \left( (1 - y_i \left( \sum_{j=1}^P z_{ij} \beta_j \right))_+ + \lambda \sum_{j=1}^P \beta_j^2 \right)$$

$$y_i \in \{1, -1\}$$



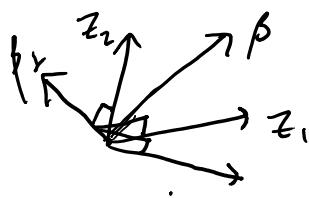
Use other losses ...

$$(\star) \min_{\beta} \frac{1}{n} \sum_{i=1}^n l(y_i, z_i^\top \beta) + \lambda \sum_{j=1}^P \beta_j^2$$

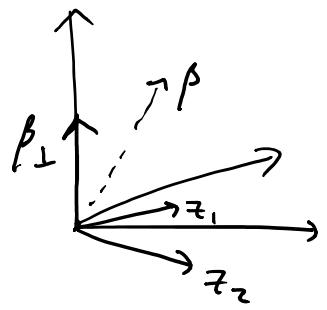
Kernel trick

$$\beta = \sum_i \alpha_i z_i + \beta_\perp$$

$$\alpha_1, \dots, \alpha_n \in \mathbb{R} \quad \beta_\perp^\top z_i = 0$$



$$z_i^\top \beta = \sum_k \alpha_k z_i^\top z_k + \cancel{z_i^\top \beta_\perp}$$



$$z_i^\top \beta = \sum_k \alpha_k z_i^\top z_k + \cancel{z_i^\top \beta} \xrightarrow{\leftarrow} 0$$

$$(\star) \min_{\alpha, \beta_\perp} \frac{1}{n} \sum_{i=1}^n l(y_i, \sum_k \alpha_k z_i^\top z_k) + \lambda \sum_{j=1}^p \underbrace{\left( \sum_k \alpha_k z_{kj} + \beta_{\perp j} \right)^2}_{(R)}$$

$$(R) = \lambda \sum_j \left( \sum_k \alpha_k z_{kj} \right)^2 + \underbrace{2 \beta_{\perp j} \left( \sum_k \alpha_k z_{kj} \right)}_{(C)} + \beta_{\perp j}^2$$

$$(C) = 2\lambda \sum_j \sum_k \alpha_k \beta_{\perp j} z_{kj} = 2\lambda \sum_k \alpha_k \underbrace{\sum_j \beta_{\perp j} z_{kj}}_{\beta_{\perp}^\top z_k = 0}$$

(dependent on  $\alpha$ )

$$(\star) \min_{\alpha, \beta_\perp} \dots + \lambda \sum_j \beta_{\perp j}^2 \quad \beta_{\perp} = 0 \quad \text{for any minimizer}$$

$$(\star) \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n l(y_i, \sum_k \alpha_k z_i^\top z_k) + \lambda \sum_j \left( \sum_k \alpha_k z_{kj} \right)^2$$

$$z_i^\top z_k = \Phi(x_i)^\top \Phi(x_k) =: k(x_i, x_k)$$

$$\sum_j \sum_{k_1, k_2} \alpha_{k_1} \alpha_{k_2} z_{kj} z_{k_2 j}$$

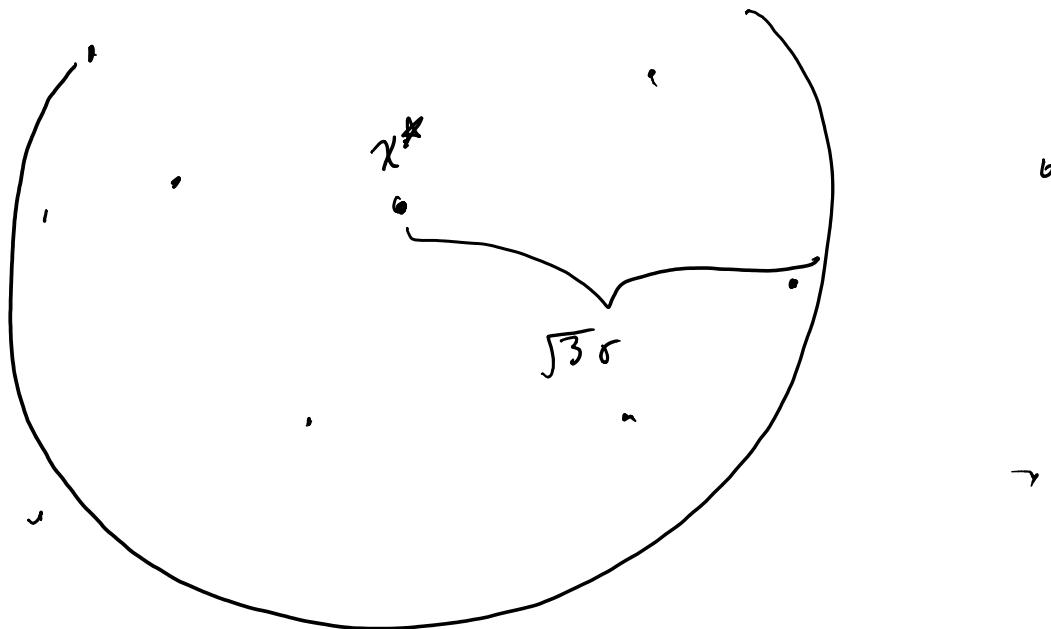
$$= \sum_{k_1} \sum_{k_2} \alpha_{k_1} \alpha_{k_2} z_{k_1}^\top z_{k_2}$$

$$(\star) \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n l(y_i, \sum_k \alpha_k k(x_i, x_k)) + \lambda \sum_{k_1, k_2} \alpha_{k_1} \alpha_{k_2} k(x_{k_1}, x_{k_2})$$

ex  $d^{\text{th}}$  degree poly kernel:  $K(x_i, x_k) = (1 + x_i^T x_k)^d$

$$\begin{aligned} d=2: \quad & (1 + x_{i1} x_{k1} + x_{i2} x_{k2})^2 = 1 + 2x_{i1} x_{k1} + 2x_{i2} x_{k2} + \\ & x_{i1}^2 x_{k1}^2 + x_{i2}^2 x_{k2}^2 + 2x_{i1} x_{i2} x_{k1} x_{k2} \\ & = (1, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}, \sqrt{2}x_{i1}x_{i2}, x_{i1}^2, x_{i2}^2)^T \\ & (\dots x_k \dots) \end{aligned}$$

ex  $K(x_i, x_k) = e^{-\|x_i - x_k\|_2^2 / \sigma^2}$  Radial basis function



$\sigma$  large: high bias low variance

$\sigma$  small: low bias high variance