

# Upcoming lecture schedule:

## Advanced topics and hands-on practices

---

- Jun. 19 (today): Classification performance measure & non-linear models [Kashima]
  - Jun. 26: **Hands-on practice** on regression and classification [Takeuchi]
  - Jul. 3: Cont'd & Neural networks [Takeuchi]
  - Jun. 10: Graph neural networks [Yamada]
  - Jun. 24: **Hands-on practice** on neural networks and graph NN [Takeuchi]
- \* It is preferable (but not mandatory) that you bring your own PC with an Internet connection for the **hands-on practices**.

# Statistical Learning Theory - Nonlinear Models -

Hisashi Kashima

## Contents:

# Classification performance measures & nonlinear models

---

- Performance measures for classification:
  - Precision / Recall
  - AUC
- Nonlinear models:
  - Simple nonlinear transformation / cross-terms
  - Kernel methods: kernel regression

# Performance Measures for Classification

# Various performance measures of classifiers:

## Should be chosen according to applications

---

- Evaluation measures to quantify the performance of a trained model especially in supervised classification
  - Accuracy, precision/recall, AUC, DCG@ $k$ , ...
- They should be appropriately chosen depending on applications
  - Classification *using* decision thresholds: accuracy, precision/recall, ...
  - Classification *without* decision thresholds: AUC, ...
  - Ranking: DCG@ $k$ , ...

## Confusion matrix:

Set of predictions on a dataset gives a confusion matrix

- A classifier makes positive (+1) or negative (−1) predictions
  - Linear classifier:  $y = \text{sign}(f(\mathbf{x}))$ ,  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
  - The larger  $f(\mathbf{x})$  is, the more strongly the classifier believes that  $\mathbf{x}$  belongs to class is  $y = +1$
- Once we have a set of predictions on a dataset, we have a confusion matrix:

|            |          | predicted label   |                   |
|------------|----------|-------------------|-------------------|
|            |          | positive          | negative          |
| true label | positive | #true positives 😊 | #false negatives  |
|            | negative | #false positives  | #true negatives 😊 |

# Basic performance measures :

## Accuracy, precision, recall

|            |          | predicted label   |                   |
|------------|----------|-------------------|-------------------|
|            |          | positive          | negative          |
| true label | positive | #true positives 😊 | #false negatives  |
|            | negative | #false positives  | #true negatives 😊 |

- Accuracy: percentage of  $\frac{\text{\#true positives} + \text{\#true negatives}}{\text{\#all predictions}}$

—In other words, averaged 0-1 loss

- Precision/Recall

—Precision =  $\frac{\text{\#true positives}}{\text{\#true positives} + \text{\#false positives}}$

—Recall =  $\frac{\text{\#true positives}}{\text{\#true positives} + \text{\#false negatives}}$

—F—measure =  $\frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

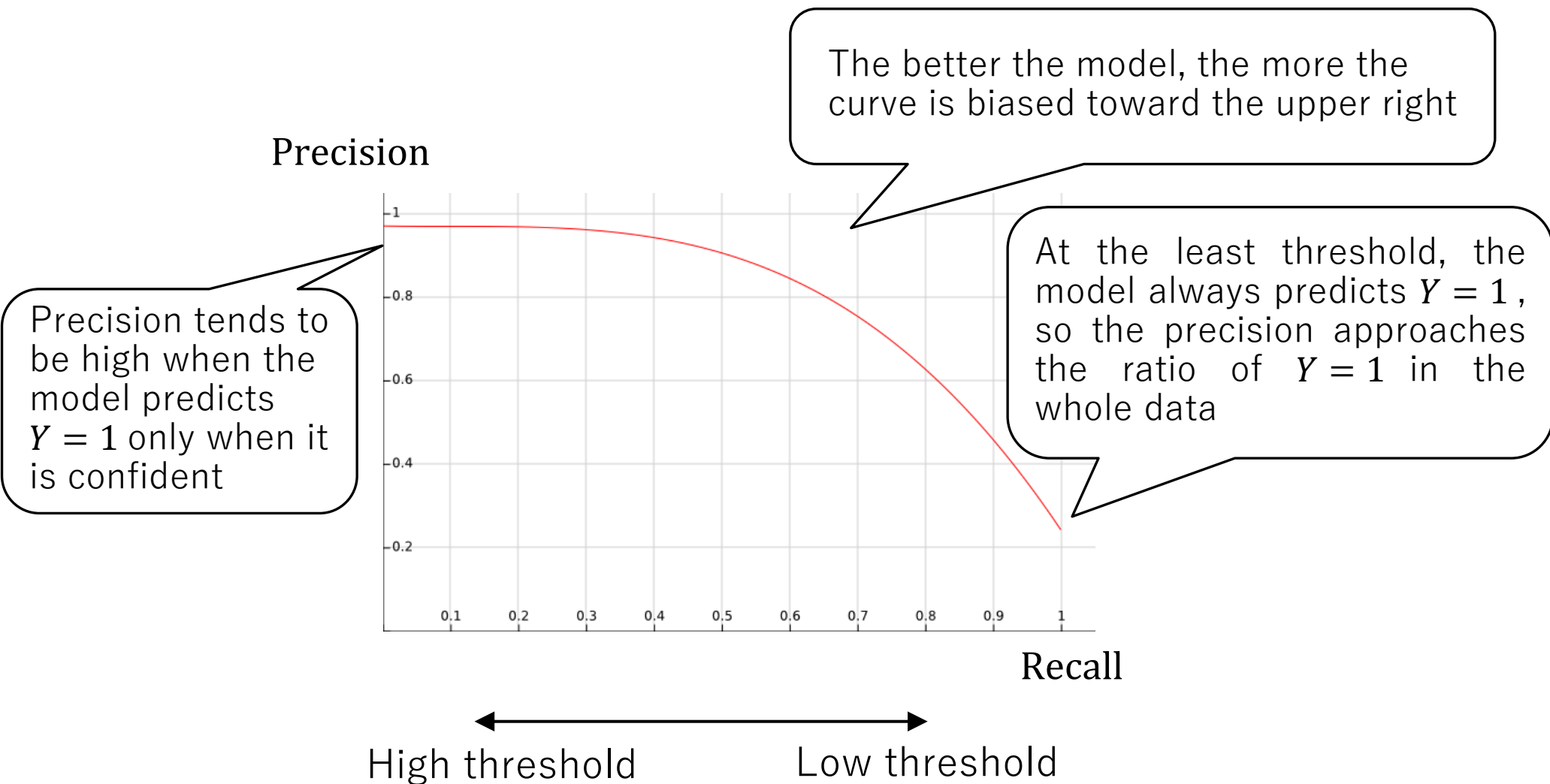
- Harmonic mean of precision and recall

Wherever he goes, there's always a murder.



Wherever there's a murder, he's always there.

# Precision-recall curve: View changes in precision & recall with different thresholds





# ROC-curve: View changes in true-positives & false-positives with different thresholds

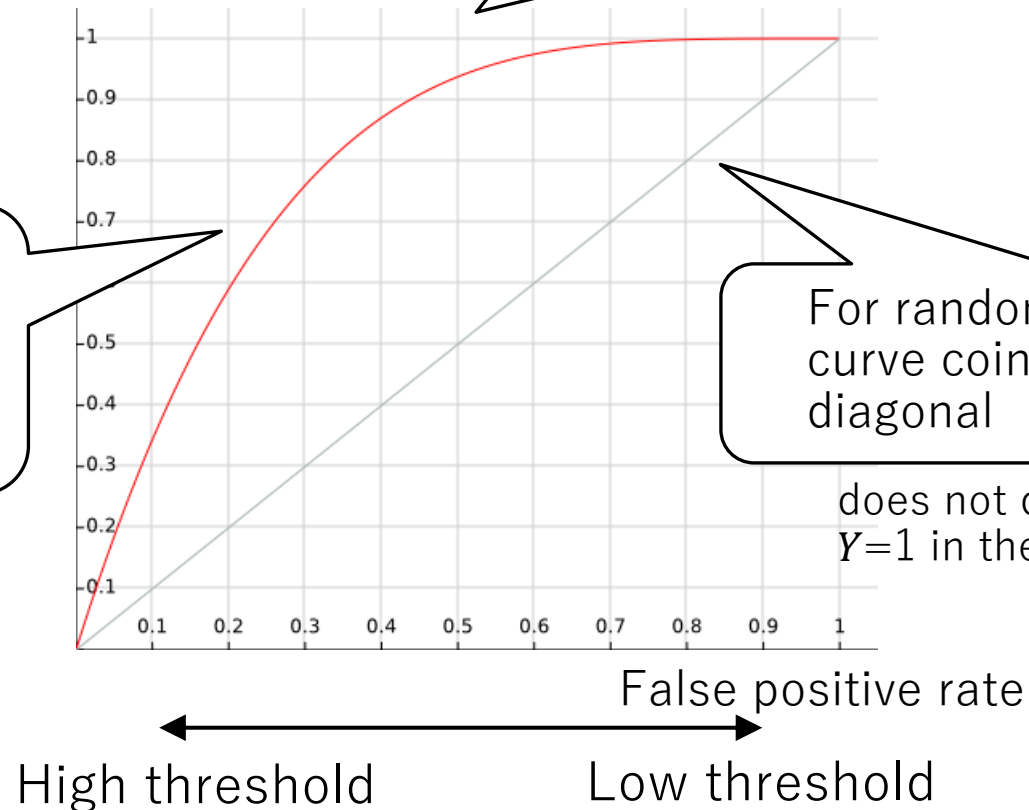
True positive rate  
(= recall)

Lowering the threshold increases the number of true positives, but also increases false positives

The better the model, the more the curve is biased toward the upper left

For random projections, the curve coincides with the diagonal

does not depend on the ratio of  $Y=1$  in the whole data

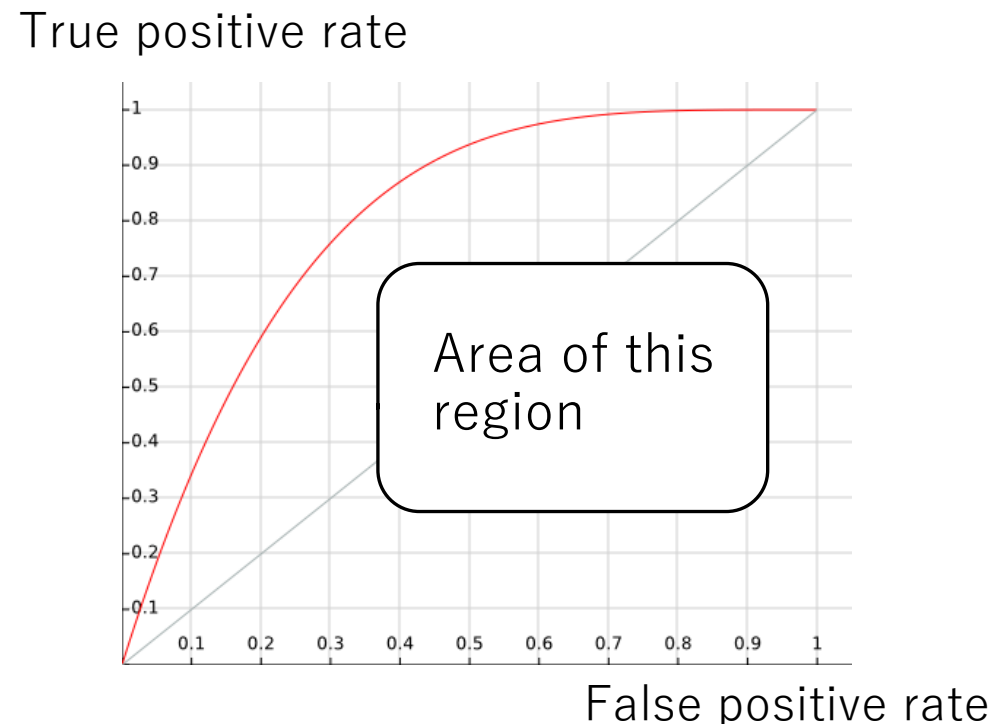
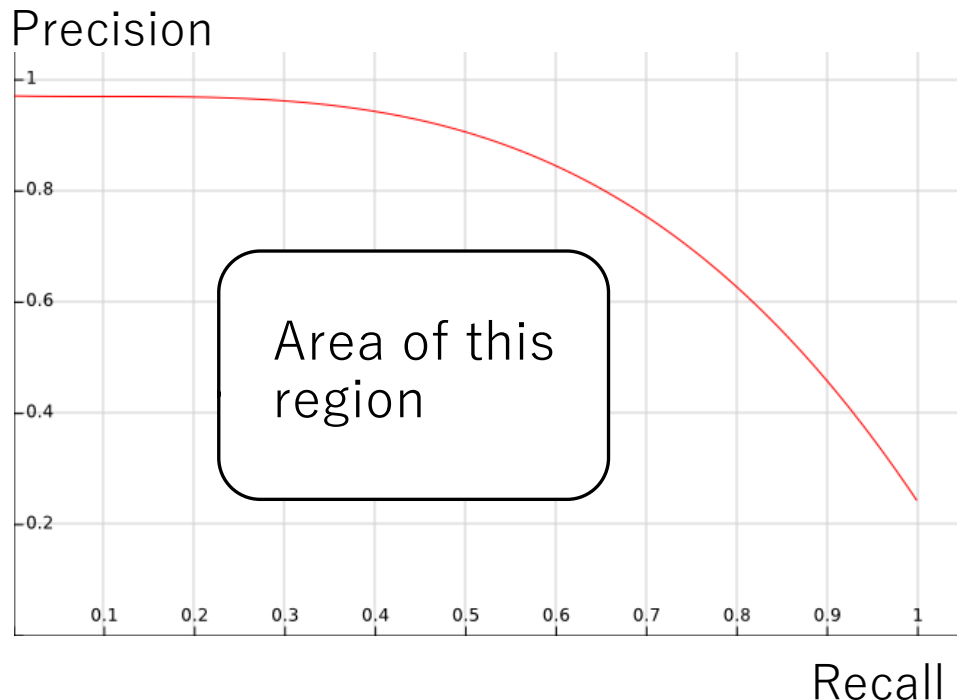


# Area under the curve:

## Performance measures independent of thresholds

- The area under the PR-curve (PR-AUC)
- The area under the ROC-curve (ROC-AUC)
  - ROC-AUC is not affected by class balance

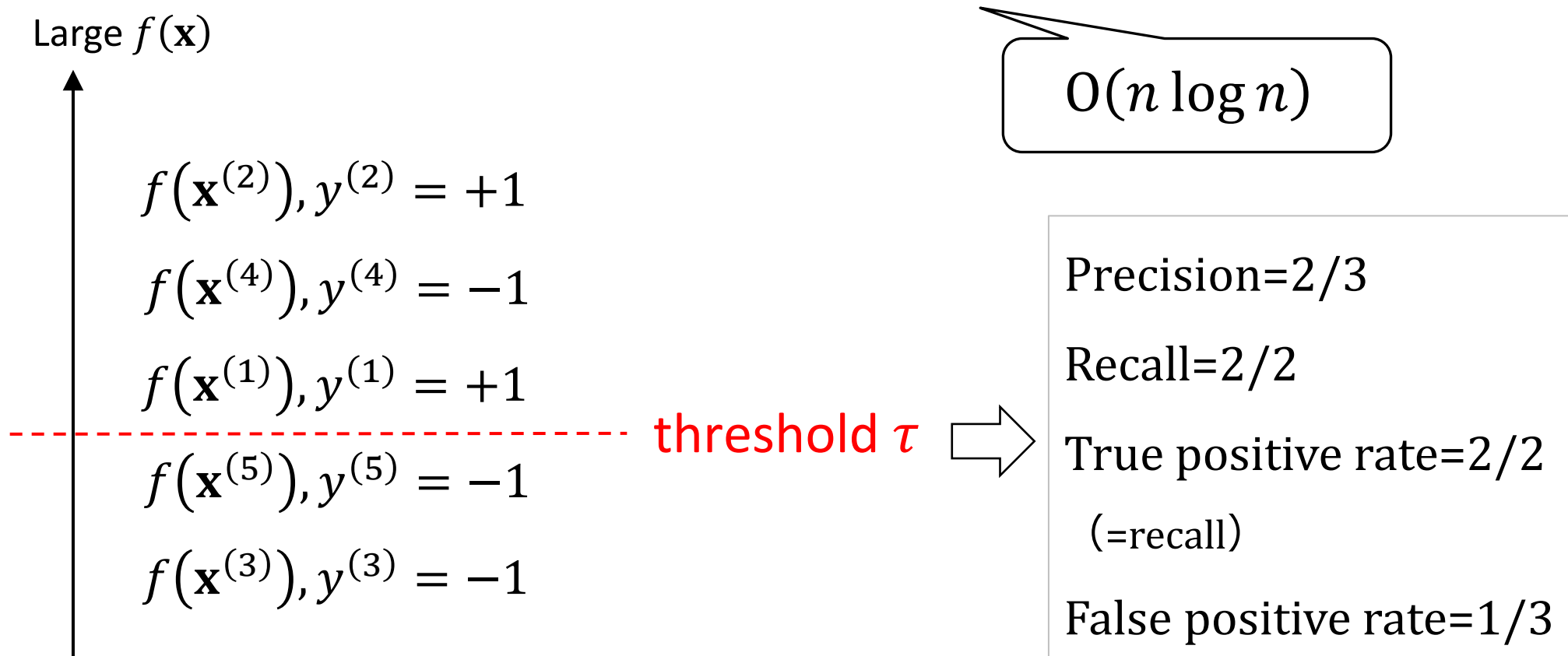
When we simply say “AUC”, we usually mean this



# Computational complexity of the performance measures:

## Sorting the model predictions

- Complexity of drawing {PR, ROC}-{curve, AUC} is equivalent to that of sorting the prediction scores  $f(\mathbf{x})$  in descending order



## Another implication of ROC-AUC:

### ROC-AUC measures ordering correctness

- ROC-AUC: Proportion of  $(i, j)$  pairs satisfying  $y^{(i)} = +1, y^{(j)} = -1$ , and  $f(\mathbf{x}^{(i)}) > f(\mathbf{x}^{(j)})$
- It checks that the test data are ranked in the correct order by  $f$ 
  - AUC=1: Perfect ranking
  - AUC=0.5: Completely random ranking
  - AUC=0: Perfectly reversed ranking
- Example: AUC=5/6
  - Among  $2 \times 3 = 6$  pos-neg pairs
  - 5 pairs are in correct order

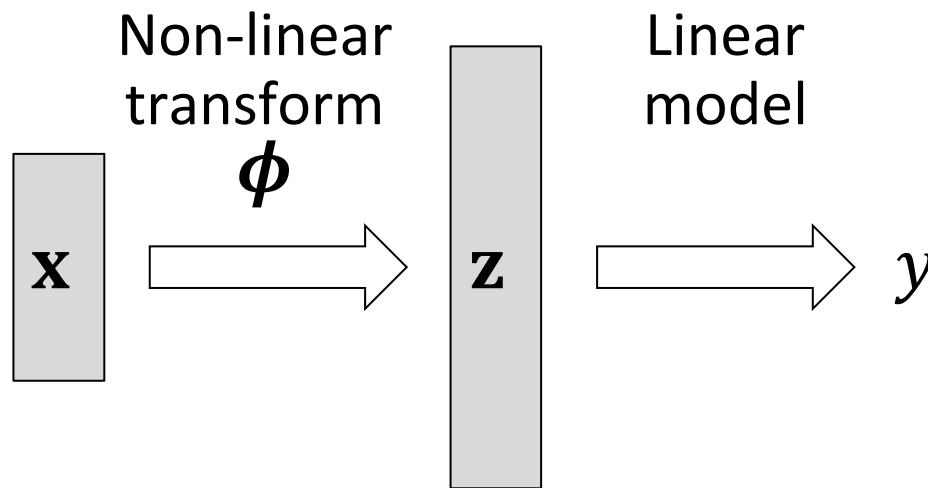
$$\begin{array}{l} f(\mathbf{x}) \uparrow \\ f(\mathbf{x}^{(2)}), y^{(2)} = +1 \\ f(\mathbf{x}^{(4)}), y^{(4)} = -1 \\ f(\mathbf{x}^{(1)}), y^{(1)} = +1 \\ f(\mathbf{x}^{(5)}), y^{(5)} = -1 \\ f(\mathbf{x}^{(3)}), y^{(3)} = -1 \end{array}$$

# Nonlinear Models

## Transformation-based nonlinear models:

Apply nonlinear transform before applying linear model

- Input vector  $\mathbf{x} \in \mathbb{R}^D$  is transformed to a new vector  $\mathbf{z} \in \mathbb{R}^{D'}$  using some nonlinear transformation function  $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$
- Linear model is applied to  $\mathbf{z}$ :
  - Linear regression  $y = \mathbf{v}^\top \mathbf{z}$ , logistic regression  $y = \sigma(\mathbf{v}^\top \mathbf{z})$



# Nonlinear regression:

## Introducing nonlinearity in linear models

---

- So far we have considered only linear models:
  - Linear regression  $y = \mathbf{w}^\top \mathbf{x}$ , logistic regression  $y = \sigma(\mathbf{w}^\top \mathbf{x})$
- How to introduce non-linearity in the models?
  1. Use of inherently nonlinear models:
    - Decision/regression tree, random forest, boosting trees
  2. Transformation-based approaches:
    - Nonlinear feature transformation
    - Kernel methods
    - Neural networks

# Nonlinear transformation of features:

## Simplest way to introduce nonlinearity in linear models

- Apply non-linear transformation:

$$- \mathbf{x} = (x) \Rightarrow \mathbf{z} = \left( \log x, e^x, x^2, \frac{1}{x}, \dots \right)^\top$$

$$- y = wx \Rightarrow y = w_1 \log x + w_2 e^x + w_3 x^2 + w_4 \frac{1}{x} + \dots$$

- It is up to the user to decide which transformations to use.



# Cross terms & factorization machine:

Can include synergetic effects among different features

- Use cross terms products  $\{x_d x_{d'}\}_{d,d'}$  of  $x_1, x_2, \dots, x_D$
- Model has a matrix parameter  $\mathbf{W}$ :

$$y = \text{Trace} \left( \begin{bmatrix} w_{1,1} & \cdots & w_{1,D} \\ \vdots & \ddots & \vdots \\ w_{D,1} & \cdots & w_{D,D} \end{bmatrix}^\top \begin{bmatrix} x_1^2 & x_1 x_2 & \cdots & x_1 x_D \\ x_2 x_1 & x_2^2 & \cdots & x_2 x_D \\ \vdots & \vdots & \ddots & \vdots \\ x_D x_1 & x_D x_2 & \cdots & x_D^2 \end{bmatrix} \right) = \mathbf{x}^\top \mathbf{W}^\top \mathbf{x}$$

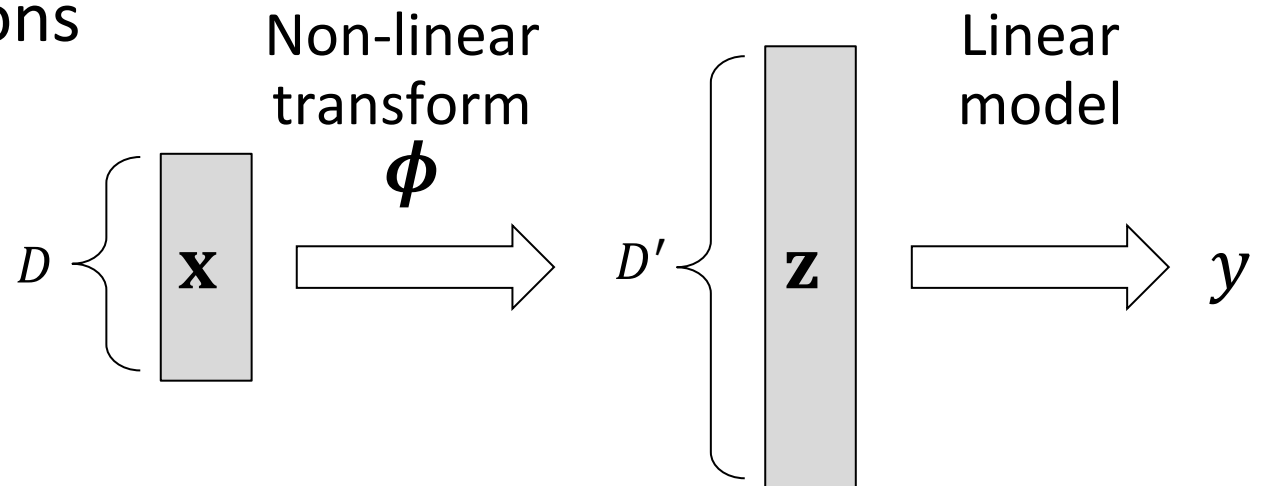
- Loss function:  $L(\mathbf{W}) = \sum_{i=1}^N \left( y^{(i)} - \mathbf{x}^{(i)\top} \mathbf{W}^\top \mathbf{x}^{(i)} \right)^2$
- Assuming low rankness of  $\mathbf{W} = \mathbf{U}\mathbf{U}^\top$  leads to factorization machine ( $O(DK)$  parameters instead of  $O(D^2)$ )

# Kernel Methods

# Kernels:

## Linear model in a high-dimensional feature space

- Kernel method is a general framework to convert a linear machine to non-linear machine
- High dimensional non-linear mapping:  $\mathbf{x} \rightarrow \mathbf{z} = \boldsymbol{\phi}(\mathbf{x})$
- Consider a linear model  $y = \mathbf{v}^\top \boldsymbol{\phi}(\mathbf{x})$  in the high dim. space
- Resolves computational difficulties caused by high dimensionality through kernel functions



# “Dual form” of linear regression model: Representation using only inner products of input vectors

- Let us construct a “kernel version” of linear regression

- Linear regression:  $y = \mathbf{w}^\top \mathbf{x}$

$D$  dimensional model

– Training data:  $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$

– Objective function:  $L(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$

– Solution:  $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

Dealing with  $D \times D$  matrix

- The computational costs are governed by  $D$

# “Dual form” of linear regression model: Representation using only inner products of input vectors

- Now we assume  $\mathbf{w} = \sum_{i=1}^N \alpha_i \mathbf{x}^{(j)}$  (weighted sum of inputs)

- (For the time being, we accept this without reason)

- $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_N)^\top$ : a new  $N$ -dimensional parameter

- We have “kernel ridge regression”:

- Model:  $y = \sum_{i=1}^N \alpha_i \langle \mathbf{x}^{(j)}, \mathbf{x} \rangle$

$N$  dimensional model

- Objective function:  $L(\mathbf{w}) = \|\mathbf{y} - \mathbf{K}\boldsymbol{\alpha}\|_2^2 + \boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha}$

- Solution:  $\mathbf{w}^* = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$

Dealing with  $N \times N$  matrix

- $\mathbf{K} = [\mathbf{K}_{i,j}] = [\langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle]$  (Kernel matrix)

## Advantage of kernel methods: Computational costs depending on the number of training data

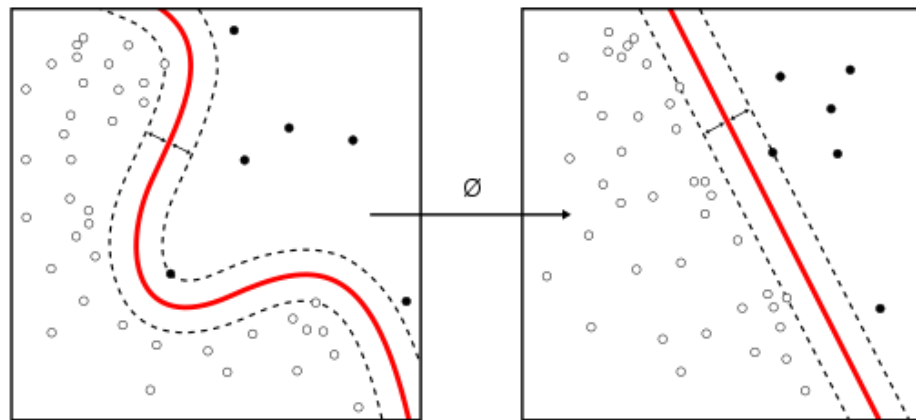
---

- What is nice about the kernel ridge regression?
  - Model/problem size depend on the size of the training data  $N$  instead of the number of dimensions  $D$
  - Computational advantage when  $D > N$
- Kernel machines access data only through kernel functions (= inner products between data)

# Kernel functions:

## Introducing non-linearity in linear models

- Consider a (nonlinear) mapping  $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$ 
  - $D$ -dimensional space to  $D' (\gg D)$ -dimensional space
  - Vector  $\mathbf{x}$  is mapped to a high-dimensional vector  $\phi(\mathbf{x})$
- Define kernel function  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \Rightarrow \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$  in the  $D'$ -dimensional space



## Advantage of kernel methods:

### Computationally efficient (when $D'$ is large)

---

- Advantage of using kernel function:

$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \boldsymbol{\phi}(\mathbf{x}^{(i)}), \boldsymbol{\phi}(\mathbf{x}^{(j)}) \rangle$$

- Usually we expect the computation cost of  $K$  depends on  $D'$ 
  - $D'$  can be high-dimensional (possibly infinite dimensional)
- If we can somehow compute  $\langle \boldsymbol{\phi}(\mathbf{x}^{(i)}), \boldsymbol{\phi}(\mathbf{x}^{(j)}) \rangle$  in time depending on  $D$ , the dimension of  $\boldsymbol{\phi}$  does not matter
- Problem size:
  - $D'$  (number of dimensions)  $\rightarrow N$  (number of data)
  - Advantageous when  $D'$  is very large or infinite



## Example of kernel functions:

### Polynomial kernel can consider high-order cross terms

- Combinatorial features: Not only the original features  $x_1, x_2, \dots, x_D$ , we use their cross terms (e.g.  $x_1 x_2$ )
  - If we consider  $M$ -th order cross terms, we have  $O(D^M)$  terms
- Polynomial kernel:  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \left( \mathbf{x}^{(i)\top} \mathbf{x}^{(j)} + c \right)^M$ 
  - E.g. when  $c = 0, M = 2, D = 2$ ,

$$\begin{aligned} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) &= \left( x_1^{(i)} x_1^{(j)} + x_2^{(i)} x_2^{(j)} \right)^2 \\ &= \left( x_1^{(i)2}, x_2^{(i)2}, \sqrt{2} x_1^{(i)} x_2^{(i)} \right) \left( x_1^{(j)2}, x_2^{(j)2}, \sqrt{2} x_1^{(j)} x_2^{(j)} \right) \end{aligned}$$

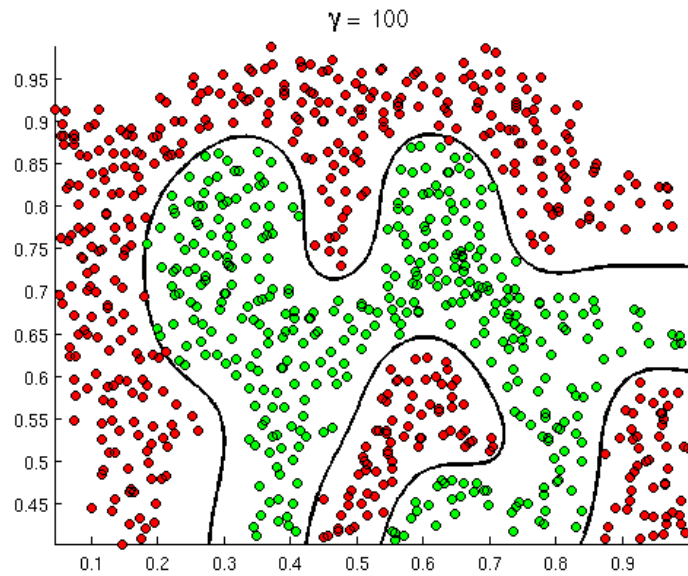
$\mathbf{x}^{(i)} = \begin{pmatrix} x_1^{(i)} \\ x_2^{(i)} \end{pmatrix}$
  - Note that it can be computed in  $O(D)$

# Example of kernel functions:

## Gaussian kernel with infinite feature space

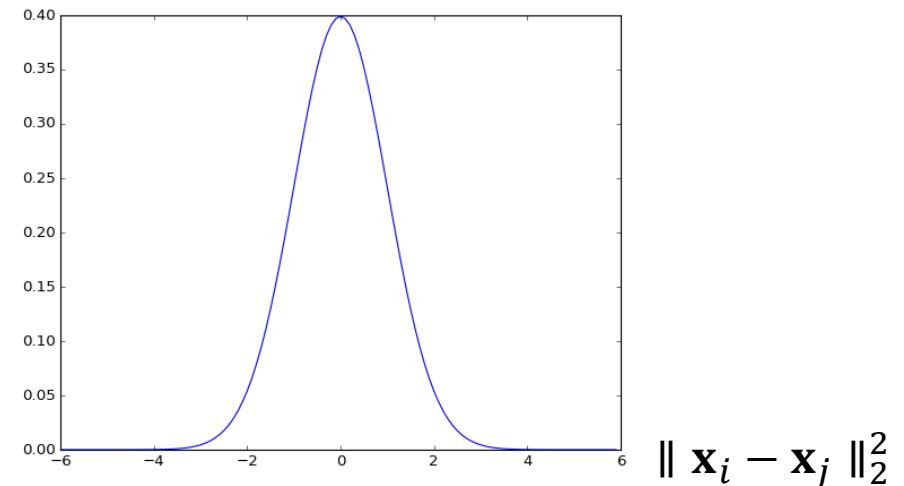
- Gaussian kernel:  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|_2^2}{\sigma}\right)$
- Can be interpreted as an inner product in an infinite-dimensional space

Discrimination surface with Gaussian kernel



<http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=MachineLearning&doc=exercises/ex8/ex8.html>

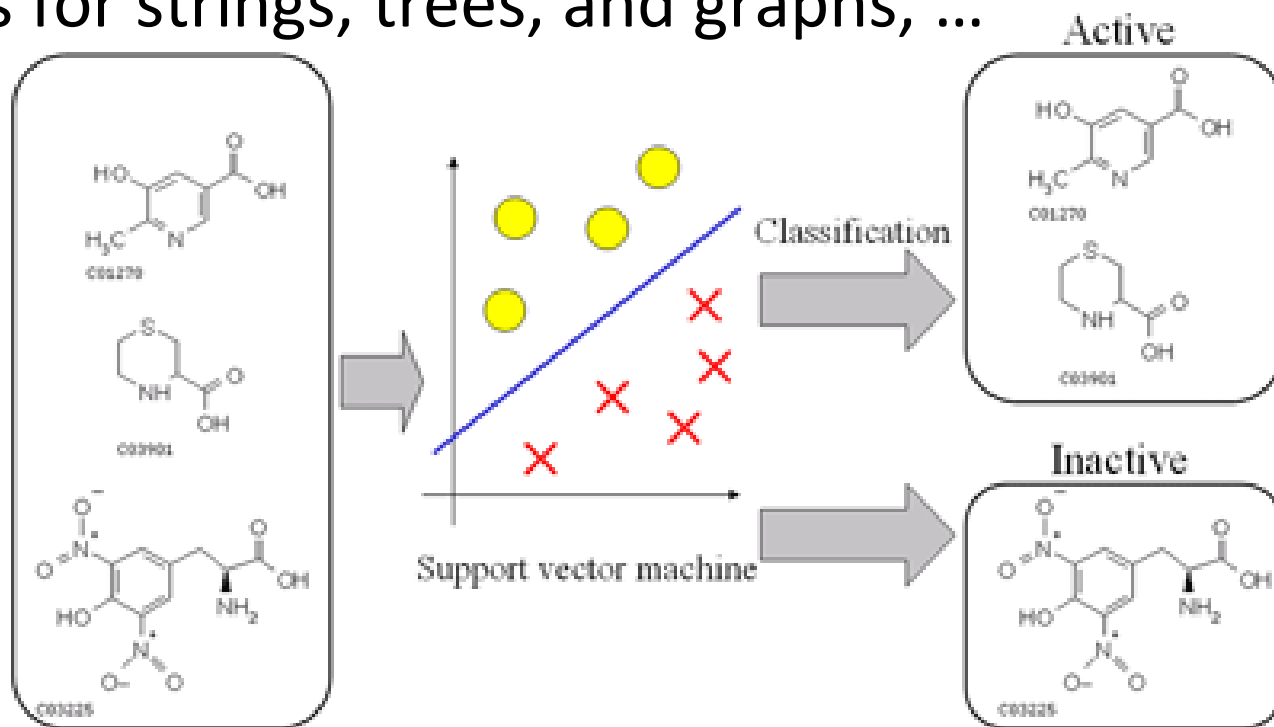
Gaussian kernel (RBF kernel)



# Kernel methods for non-vectorial data:

## Kernels for sequences, trees, and graphs

- Kernel methods can handle any kinds of objects (even non-vectorial objects) as long as efficiently computable kernel functions are available
  - Kernels for strings, trees, and graphs, ...



[http://www.bic.kyoto-u.ac.jp/coe/img/akutsu\\_fig\\_e\\_02.gif](http://www.bic.kyoto-u.ac.jp/coe/img/akutsu_fig_e_02.gif)

# Representer theorem:

## Theoretical underpinning of kernel methods

---

- Can we use some similarity function as a kernel function?
  - Yes (under *certain conditions*)
- Kernel methods rely on the fact that the optimal parameter is represented as a linear combination of input vectors:

$$\mathbf{w} = \sum_{i=1}^N \alpha_i \mathbf{x}^{(i)}$$

- Representer theorem guarantees this (if we use L2-regularizer)

## (Simple) proof of representer theorem:

Obj. func. depends only on linear combination of inputs

- Assumption: Loss  $\ell$  for  $i$ -th data depends only on  $\mathbf{w}^\top \mathbf{x}^{(i)}$ 
  - Objective function:  $L(\mathbf{w}) = \sum_{i=1}^N \ell(\mathbf{w}^\top \mathbf{x}^{(i)}) + \lambda \|\mathbf{w}\|_2^2$
- Divide the optimal parameter  $\mathbf{w}^*$  into two parts  $\mathbf{w} + \mathbf{w}^\perp$ :
  - $\mathbf{w}$ : Linear combination of input data  $\{\mathbf{x}^{(i)}\}_i$
  - $\mathbf{w}^\perp$ : Other parts (orthogonal to all input data  $\{\mathbf{x}^{(i)}\}$ )
- $L(\mathbf{w}^*)$  depends only on  $\mathbf{w}$ :  $\sum_{i=1}^N \ell(\mathbf{w}^{*\top} \mathbf{x}^{(i)}) + \lambda \|\mathbf{w}^*\|_2^2$ 
$$= \sum_{i=1}^N \ell \left( \mathbf{w}^\top \mathbf{x}^{(i)} + \underbrace{\mathbf{w}^\perp{}^\top \mathbf{x}^{(i)}}_{=0} \right) + \lambda (\underbrace{\|\mathbf{w}\|_2^2}_{=0} + \underbrace{2\mathbf{w}^\top \mathbf{w}^\perp}_{=0} + \underbrace{\|\mathbf{w}^\perp\|_2^2}_{\text{Minimized to } =0})$$