

# 7. Multi-class Classification and Kernel Methods

Dongwoo Kim

[dongwoo.kim@postech.ac.kr](mailto:dongwoo.kim@postech.ac.kr)

CSED515 - 2023 Spring

# Application of Classification

## Binary classification

- ▶ Image classification: Is the image showing a cat? Yes/No
- ▶ Edge detection: Is a given pixel at a boundary? Yes/No
- ▶ Object detection: Is there a cat in the image? Yes/No
- ▶ ...

## Multi-class classification

- ▶ Which object is illustrated? Cat/Dog/Tiger/Zebra...



# Table of Contents

- 1 Understanding multi-class logistic regression
- 2 Getting to know multi-class SVM
- 3 Learning the kernel trick
- 4 Noise-Free Gaussian Process Regression

## Recall) Discriminant Function

A description of classifier:

assign  $x$  to  $\mathcal{C}_k$  if  $f_k(x) > f_j(x)$  for all  $j \neq k$

- ▶  $\{f_k\}_{k=1,\dots,K}$ : discriminant function/score
  - ▶ e.g.,  $f_k(x) = p(\mathcal{C}_k | x)$
- ▶ i.e., classifier  $h(x) = \arg \max_k f_k(x)$
- ▶ If  $f_k(\cdot)$ 's are a set of discriminant functions, then for any monotonically increasing function  $g : \mathbb{R} \rightarrow \mathbb{R}$   $g_k(\cdot) = g(f_k(\cdot))$ 's are also a set of discriminant functions
  - ▶ e.g.1,  $a f_k(x) + b$  for any  $a > 0$  and  $b$
  - ▶ e.g.2,  $\log f_k(\cdot)$
  - ▶ ...

## Recall) Discriminant Function: Bayes Decision Rule

Bayes decision rule:

$$f_k(x) = p(\mathcal{C}_k | x) \text{ or equivalently, } f_k(x) = p(x | \mathcal{C}_k)p(\mathcal{C}_k)$$

- Classify  $x \in \mathcal{C}_k$  if  $p(\mathcal{C}_k | x) > p(\mathcal{C}_j | x)$  for all  $j \neq k$ .
- We can use discriminant function  $f_k(x) = p(x | \mathcal{C}_k)p(\mathcal{C}_k)$  without **normalization factor** thanks to Bayes rule:

$$p(\mathcal{C}_k | x) = \frac{\overbrace{p(x | \mathcal{C}_k)}^{\text{Class-conditional density}} \overbrace{p(\mathcal{C}_k)}^{\text{Prior}}}{\underbrace{\sum_{j=1}^K p(x | \mathcal{C}_j)p(\mathcal{C}_j)}_{\text{Normalization factor}}}$$

In general, a classification task can be understand as a task of learning  $p(y = k | x)$  (or score function) for each class  $k$

# Multi-class Logistic Regression for Classification

- ▶ A model of multinomial distribution over  $K$ -class  $y \in \{1, 2, \dots, K\}$ :
  - ▶ Using  $K$  weight vectors, denoted by  $\mathbf{w}_{(y)}$ 's, for every class

$$p(y = k|x) = \frac{\exp(\mathbf{w}_{(k)}^\top \phi(x))}{\sum_{j=1}^K \exp(\mathbf{w}_{(j)}^\top \phi(x))}$$

- ▶ Classifier (Bayes decision rule):

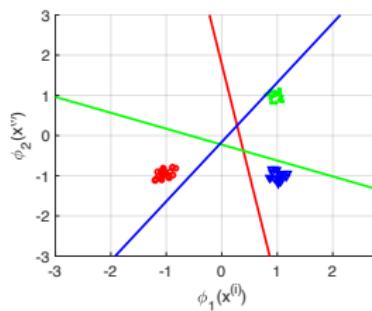
$$\arg \max_k p(y = k|x) = \arg \max_k \mathbf{w}_{(k)}^\top \phi(x)$$

- ▶ Learning  $\mathbf{w} = (\mathbf{w}_{(y)})_{y=1,\dots,K}$  is maximizing the likelihood (assuming i.i.d. data points) as before:

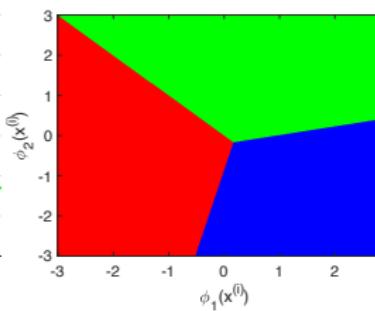
$$\arg \max_{\mathbf{w}} \prod_{i=1,\dots,N} p(y = y_i | x^{(i)}) = \arg \min_{\mathbf{w}} -\log p(y = y_i | x^{(i)})$$

# An Example of Multi-class Logistic Regression

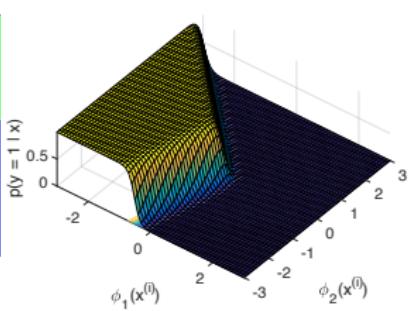
Setup



Decision Boundary



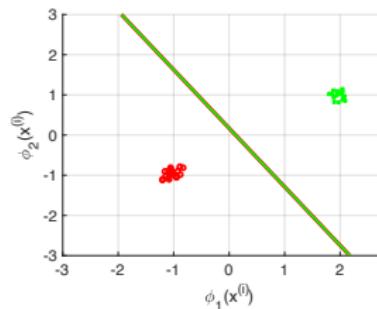
Probability for  $k = 1$



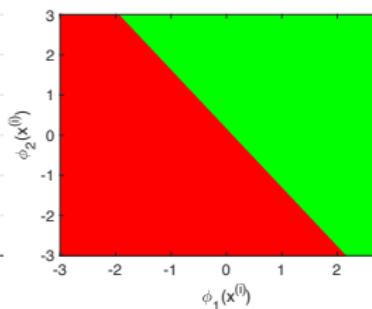
# An Example of Multi-class Logistic Regression

What if  $K = 2$ ?

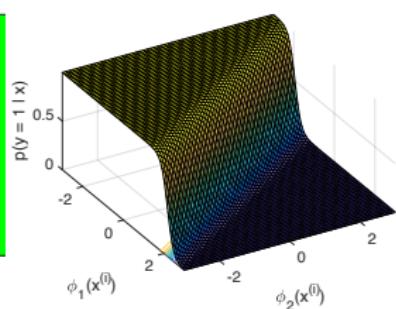
Setup



Decision Boundary



Probability for  $k = 1$



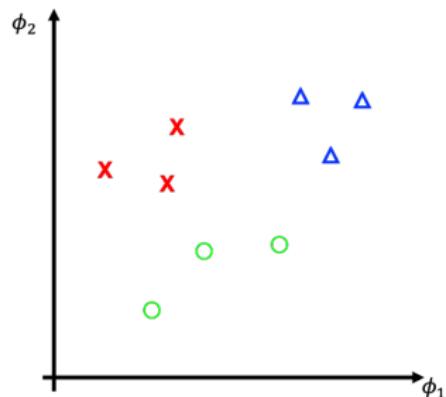
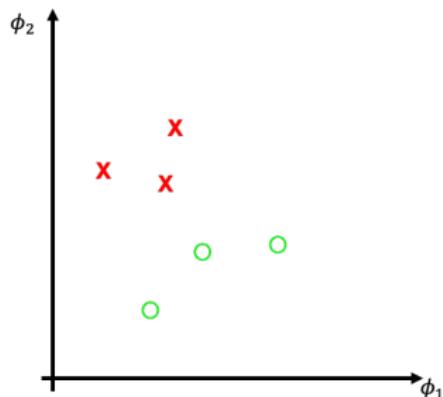
We have  $\mathbf{w}_{(1)} = -\mathbf{w}_{(2)}$

# Table of Contents

- 1 Understanding multi-class logistic regression
- 2 Getting to know multi-class SVM
- 3 Learning the kernel trick
- 4 Noise-Free Gaussian Process Regression

## Binary SVM: Drawing Hyperplane

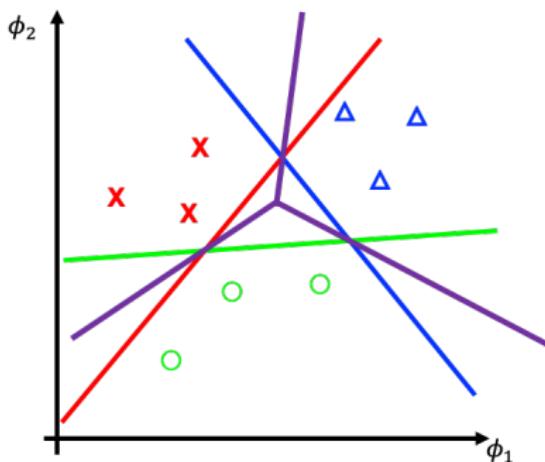
Combination of binary SVMs =  $K$ -class SVM?



## Approach1: One vs. Rest (OvR)

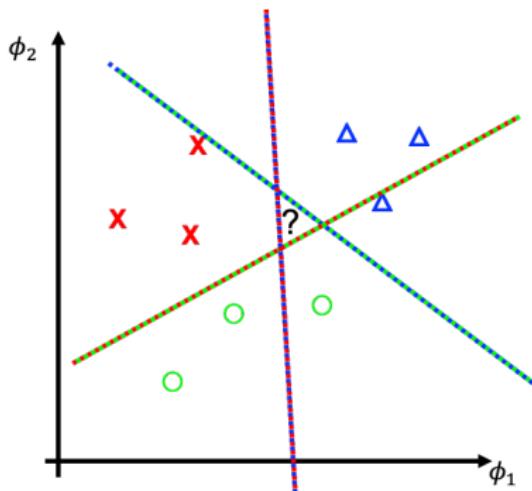
One-vs-Rest, OvR; One-vs-All, OvA; or One-against-All, OAA

- ▶ Learn  $K$  SVM's, each of which is denoted by  $(\mathbf{w}_{(k)}, b_{(k)})$  for binary classification between class  $k$  (positive) and others (negative)
- ▶ For new data  $\mathbf{x}$ , output  $\arg \max_{k \in [K]} (\mathbf{w}_{(k)}^\top \phi(\mathbf{x}) + b_{(k)})$



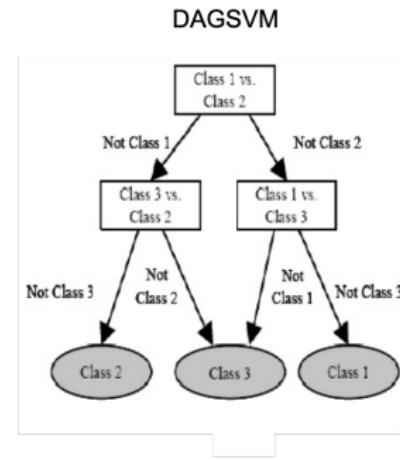
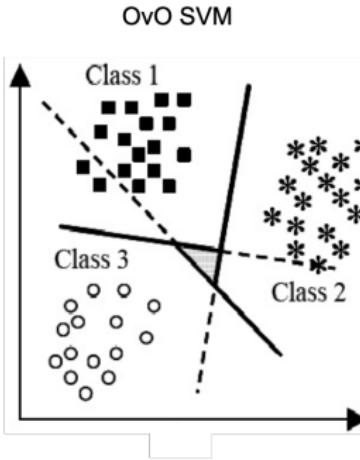
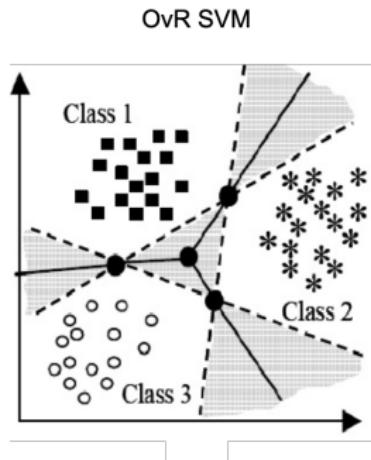
## Approach2: One vs. (another) One (OvO)

- ▶ Learn  $K(K - 1)/2$  SVM's, each of which is denoted by  $(\mathbf{w}_{(jk)}, b_{(jk)})$  binary classification between two different classes  $j$  and  $k$ , distinguishing classes  $k$  and  $j$
- ▶ For new data  $\mathbf{x}$ , take the majority vote of  $K(K - 1)/2$  SVM's with some tie-breaking rule



# OvR and OvO SVM

- ▶ OvR SVM has class **imbalance problems** as a class is much smaller than the others.
- ▶ Each of OvO SVM learns a part of dataset and **computationally expensive**.



## Approach3 by Weston & Watkins 1999 (WW-SVM)

- ▶ Solve

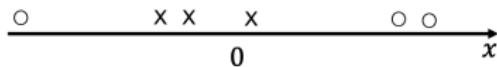
$$\begin{aligned} \min_{\mathbf{w}, \mathbf{b}, \zeta \geq 0} \quad & \frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_{(k)}\|_2^2 + C \sum_{i=1}^N \sum_{k \neq y_i} \zeta_{(k)}^{(i)} \\ \text{s.t.} \quad & \mathbf{w}_{(y_i)}^\top \phi(\mathbf{x}_i) + b_{(y_i)} - (\mathbf{w}_{(k)}^\top \phi(\mathbf{x}_i) + b_{(k)}) \geq 2 - \zeta_{(k)}^{(i)} \end{aligned}$$

- ▶ For new data  $\mathbf{x}$ , output  $\arg \max_{k \in [K]} (\mathbf{w}_{(k)}^\top \phi(\mathbf{x}) + b_{(k)})$
- ▶ This may resolve the imbalance issue in OvR SVM, and the local-view issue in OvO SVM
- ▶ What are the dual problem, and the support vectors in WW SVM?  
c.f.,

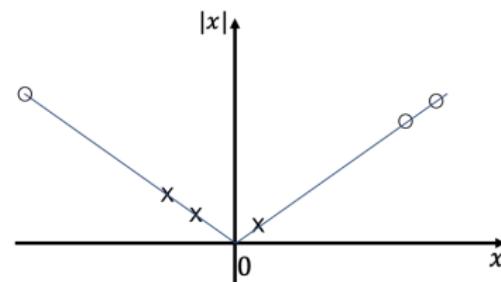
<https://www.csie.ntu.edu.tw/~cjlin/papers/multisvm.pdf>

# Are We Done with SVM?

- ▶ Still, the representation power of SVM is limited to hyperplanes
  - ▶ Option1. allow slack ( $\zeta^{(i)}$ 's) a.k.a. soft-margin SVM for dataset which are not linearly separable
  - ▶ Option2. extract more/other features  $\phi$  or use kernel trick



Not linearly separable



Linearly separable

# Table of Contents

- 1 Understanding multi-class logistic regression
- 2 Getting to know multi-class SVM
- 3 Learning the kernel trick
- 4 Noise-Free Gaussian Process Regression

## Dual Problem of Binary Soft-Margin SVM - L1 Norm

Dual of soft-margin SVM with  $\ell_1$  norm:

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) + \sum_{i=1}^N \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \quad \forall i = 1, \dots, N \end{aligned}$$

- ▶ Note that the problem depends on  $\phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$ 's
- ▶ Further more, the solution also depends on  $\phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$ 's
- ▶ Kernel trick: use  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  instead of  $\phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$

## Kernel Trick

- ▶ Computing exact feature  $\phi(\mathbf{x})$  for the entire dataset could be expensive.
  - ▶ i.e.  $\phi(\mathbf{x}) = [1, x_1, x_2, \dots, x_1^2, x_1x_2, x_1x_3, \dots, x_1^3, x_1^2x_2, \dots]^T$
  - ▶ The dimension of  $\mathbf{w}$  increases as we expand the feature map.
- ▶ Kernel trick: replace the inner product  $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$  by kernel function,  $\kappa(\mathbf{x}, \mathbf{x}')$
- ▶ We don't need to compute exact feature map  $\phi(\mathbf{x})$  given kernel function.
  - ▶ The number of Lagrangian multipliers  $\alpha$  doesn't depend on the dimension of feature map.

## Dual Objective Revised

- ▶ With kernel function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^N \alpha_i$$

$$\text{s.t. } \sum_{i=1}^N y_i \alpha_i = 0$$

$$0 \leq \alpha_i \leq C \quad \text{for all } i = 1, \dots, N$$

- ▶ We don't need to evaluate feature map during training.

## Kernels - Least Square with Gradient Descent

Kernel method can be applied to any ML algorithms that utilize the inner product between inputs.

For example, a gradient descent step for the least square method is

$$\theta := \theta + \eta \sum_{n=1}^N (y_n - \theta^\top \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)$$

Note that<sup>1</sup>  $\theta$  is a linear combination of  $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)$  if the initial value of  $\theta = 0$ .

$$\theta = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$$

---

<sup>1</sup>We can prove this inductively.

## Kernels - Least Square

The update step can be rewritten as

$$\begin{aligned}\theta &:= \theta + \eta \sum_{n=1}^N (y_n - \theta^\top \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) \\ &= \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n) + \eta \sum_{n=1}^N (y_n - \theta^\top \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) \\ &= \sum_{n=1}^N (\alpha_n + \eta(y_n - \theta^\top \phi(\mathbf{x}_n))) \phi(\mathbf{x}_n)\end{aligned}$$

New  $\alpha_n$  depends on the old one via:

$$\alpha_n := \alpha_n + \eta(y_n - \theta^\top \phi(\mathbf{x}_n))$$

## Kernels - Least Square

From

$$\alpha_n := \alpha_n + \eta(y_n - \theta^\top \phi(\mathbf{x}_n))$$

Replacing  $\theta$  with  $\theta = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$  gives

$$\alpha_i := \alpha_i + \eta \left( y_n - \sum_{j=1}^N \alpha_j \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_i) \right)$$

Result in

$$\alpha_i := \alpha_i + \eta \left( y_n - \sum_{j=1}^N \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right)$$

# Properties of Kernels

We don't need to define feature function  $\phi$  explicitly as long as we have a valid kernel function  $\kappa$ .

What are the properties of the **valid** kernel function then?

## Theorem (Mercer)

Let  $\kappa : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be given. Then for  $\kappa$  to be a valid (Mercer) kernel, it is necessary and sufficient that for any  $\mathbf{x}_1, \dots, \mathbf{x}_n, (n < \infty)$ , the corresponding kernel matrix is symmetric positive semi-definite<sup>2</sup>.

From kernel function, we can construct the kernel matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$  filled with  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ .

---

<sup>2</sup>Kernel matrices are p.s.d

## Kernels as Similarity Metrics

Intuitively, if  $\phi(\mathbf{x}_i)$  and  $\phi(\mathbf{x}_j)$  are close together, then we might expect  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  to be large.

Conversely, if  $\phi(\mathbf{x}_i)$  and  $\phi(\mathbf{x}_j)$  are nearly orthogonal to each other, then we might expect  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  to be small.

So, we can think of  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  as some measurement of how similar are  $\phi(\mathbf{x}_i)$  and  $\phi(\mathbf{x}_j)$ . For instance,

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{||\mathbf{x}_i - \mathbf{x}_j||^2}{2\sigma^2}\right)$$

This kernel is called Gaussian (Radial Basis Function: RBF) kernel.

## Example: Polynomial Kernel

- ▶ Polynomial kernel of degree  $p$ :

$$\kappa(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + 1)^p = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

- ▶ If  $\mathbf{x} = [x_1, x_2]^\top \in \mathbb{R}^2$  and  $p = 2$ , the corresponding feature map is

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^\top$$

- ▶ Instead of converting  $\mathbf{x}$  to feature  $\phi(\mathbf{x})$  and taking inner product,
- ▶ It would be much faster to compute the kernel function directly

## Example: RBF Kernel

- ▶ Radial basis function (RBF) kernel:

$$\kappa(x, x') = \exp\left(-\frac{(x - x')^2}{2\sigma^2}\right) = \phi(x)^\top \phi(x')$$

- ▶ The corresponding feature map  $\phi$  has infinite dimensions

$$\phi(x) = \exp(-x^2/2\sigma^2) \left[ 1, \sqrt{\frac{1}{1!\sigma^2}}x, \sqrt{\frac{1}{2!\sigma^4}}x^2, \sqrt{\frac{1}{3!\sigma^6}}x^3, \dots \right]$$

- ▶ We can still evaluate kernel, but cannot construct feature map.

## Additional Examples 1

Linear kernel:

$$\kappa(x^{(i)}, x^{(j)}) = x^{(i)\top} x^{(j)}$$

Squared exponential (Gaussian) kernel:

$$\kappa(x^{(i)}, x^{(j)}) = \exp\left(-\frac{1}{2}\left(x^{(i)} - x^{(j)}\right)^\top \Sigma^{-1} \left(x^{(i)} - x^{(j)}\right)\right)$$

Sigmoid kernel:

$$\kappa(x^{(i)}, x^{(j)}) = \tanh\left(a \cdot x^{(i)\top} x^{(j)} + b\right)$$

## Additional Examples 2

Positive scaling:

$$\kappa(x^{(i)}, x^{(j)}) = c \kappa_1(x^{(i)}, x^{(j)})$$

Exponentiation:

$$\kappa(x^{(i)}, x^{(j)}) = \exp(\kappa_1(x^{(i)}, x^{(j)}))$$

Addition:

$$\kappa(x^{(i)}, x^{(j)}) = \kappa_1(x^{(i)}, x^{(j)}) + \kappa_2(x^{(i)}, x^{(j)})$$

Multiplication with function:

$$\kappa(x^{(i)}, x^{(j)}) = f(x^{(i)}) \kappa_1(x^{(i)}, x^{(j)}) f(x^{(j)})$$

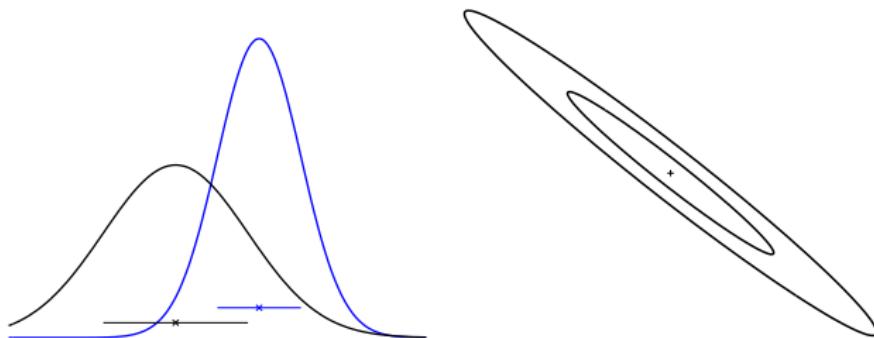
Multiplication:

$$\kappa(x^{(i)}, x^{(j)}) = \kappa_1(x^{(i)}, x^{(j)}) \kappa_2(x^{(i)}, x^{(j)})$$

# Table of Contents

- 1 Understanding multi-class logistic regression
- 2 Getting to know multi-class SVM
- 3 Learning the kernel trick
- 4 Noise-Free Gaussian Process Regression

# The Gaussian Distribution

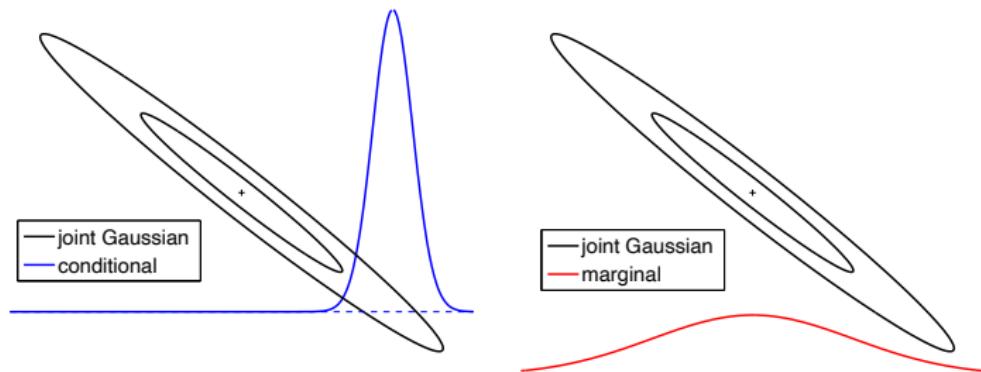


The Gaussian distribution is given by

$$p(\mathbf{x} | \mu, \Sigma) = \mathcal{N}(\mu, \Sigma) = (2\pi)^{-D/2} |\Sigma|^{-1/2} \exp \left( -\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) \right)$$

where  $\mu$  is the mean vector and  $\Sigma$  the covariance matrix.

# Conditionals and Marginals of a Gaussian



Both the conditionals and the marginals of a joint Gaussian are again Gaussian.

## Gaussian Distribution as a Distribution over Function

- ▶ Informally, function  $f : \mathcal{X} \rightarrow \mathbb{R}$  can be seen as an **infinite dimensional** vector. For example, if  $\mathcal{X} = \mathbb{R}$ , we can think

$$f(x) = [f(x=1), f(x=2), f(x=3), \dots]$$

- ▶ Given dataset  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , if the collection of outputs  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  are drawn from Gaussian distribution

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

We say the Gaussian distribution defines a distribution over function outputs.

# Gaussian Distribution as a Distribution over Function

- ▶ Informally, function  $f : \mathcal{X} \rightarrow \mathbb{R}$  can be seen as an **infinite dimensional** vector. For example, if  $\mathcal{X} = \mathbb{R}$ , we can think

$$f(x) = [f(x=1), f(x=2), f(x=3), \dots]$$

- ▶ Given dataset  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , if the collection of outputs  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  are drawn from Gaussian distribution

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

We say the Gaussian distribution defines a distribution over function outputs.

- ▶ But where are the inputs  $(x_1, x_2, \dots, x_n)$ ?
- ▶ And how can we predict output given new input  $x_*$ ?

## Kernel as Covariance Matrix

Note that if we define a Gaussian over the collection of outputs, the distribution is defined **without inputs**.

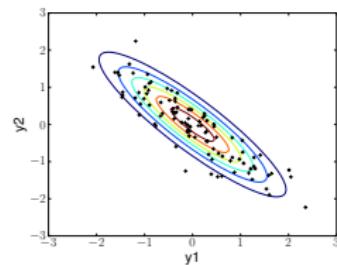
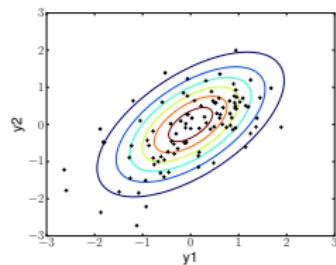
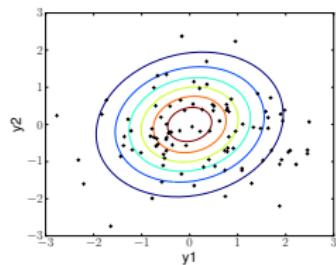
$$\mathbf{y} \sim \mathcal{N}(\mu, \Sigma)$$

- ▶ Idea: construct the covariance matrix using a kernel function  $\kappa(x, x')$ !

$$\Sigma_{i,j} = \mathbf{K}_{i,j} = \kappa(x_i, x_j)$$

- ▶ Why:
  1. **Smoothness assumption:** if two inputs  $x_1$  and  $x_2$  are close, i.e.  $x_1 \approx x_2$ , then their outputs are close, i.e.  $f(x_1) \approx f(x_2)$ .
  2. The covariance measures the correlation between two dimensions.
  3. If  $\mathbf{K}_{i,j}$  large, then  $y_i$  and  $y_j$  are likely to be similar.
  4. If  $\mathbf{K}_{i,j}$  small, then  $y_i$  and  $y_j$  are likely to be uncorrelated.

## A 2D Gaussian: Varying the Covariance



$$\mathbf{K} = \begin{bmatrix} 1 & 0.14 \\ 0.14 & 1 \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 1 \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}$$

## Putting Them Together

- Given dataset  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , we can define a probability density function over the outputs:

$$p\left(\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{bmatrix}, \begin{bmatrix} \kappa(x_1, x_1) & \dots & \kappa(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \kappa(x_n, x_1) & \dots & \kappa(x_n, x_n) \end{bmatrix}\right)$$

- Often we assume  $\mu(x) = 0, \forall x$  as a prior of function.
- We can define a distribution over functions!

## Putting Them Together

- Given dataset  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , we can define a probability density function over the outputs:

$$p\left(\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{bmatrix}, \begin{bmatrix} \kappa(x_1, x_1) & \dots & \kappa(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \kappa(x_n, x_1) & \dots & \kappa(x_n, x_n) \end{bmatrix}\right)$$

- Often we assume  $\mu(x) = 0, \forall x$  as a prior of function.
- We can define a distribution over functions!
- How can we predict output given new input  $x_*$ ?
  - Do you remember Gaussian identity? (from linear regression)

## Lemma (Gaussian identities<sup>3</sup>)

Let the joint distribution of  $z = [x^\top, y^\top]^\top$  be

$$z = \begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} a \\ b \end{bmatrix}, \begin{bmatrix} A & C \\ C^\top & B \end{bmatrix}\right) \quad \text{with cross-covariance } C = \mathbb{E}[xy^\top],$$

so that  $x \sim \mathcal{N}(a, A)$  and  $y \sim \mathcal{N}(b, B)$ . Then, the conditional distributions are given as:

$$p(x | y) = \mathcal{N}(x | a + CB^{-1}(y - b), A - CB^{-1}C^\top),$$

$$p(y | x) = \mathcal{N}(y | b + C^\top A^{-1}(x - a), B - C^\top A^{-1}C).$$

---

<sup>3</sup>c.f. Pattern Recognition and Machine Learning, Sec 2.3.1

## Prediction with Gaussian Identity

Let  $x_*$  be the test input and  $y_*$  be the predicted value. The joint distribution of  $\mathbf{y}$  and  $y_*$  given  $D$  is

$$p(\mathbf{y}, y_* \mid \mathbf{x}, \mathbf{x}^*) = \mathcal{N}\left(\begin{array}{c|c} \mathbf{y} & \mathbf{0}, \\ y_* & \left[ \begin{array}{cc} \mathbf{K}_{\mathbf{x}, \mathbf{x}} & \mathbf{K}_{\mathbf{x}, x_*} \\ \mathbf{K}_{x_*, \mathbf{x}} & \mathbf{K}_{x_*, x_*} \end{array} \right] \end{array}\right)$$

where  $\mathbf{K}_{x_*, \mathbf{x}}$  is

$$\mathbf{K}_{x_*, \mathbf{x}} = [\kappa(x_*, x_1), \dots, \kappa(x_*, x_n)]$$

## Prediction with Gaussian Identity

Let  $x_*$  be the test input and  $y_*$  be the predicted value. The joint distribution of  $\mathbf{y}$  and  $y_*$  given  $D$  is

$$p(\mathbf{y}, y_* \mid \mathbf{x}, \mathbf{x}^*) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \mid \mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{x}, \mathbf{x}} & \mathbf{K}_{\mathbf{x}, x_*} \\ \mathbf{K}_{x_*, \mathbf{x}} & \mathbf{K}_{x_*, x_*} \end{bmatrix}\right)$$

where  $\mathbf{K}_{x_*, \mathbf{x}}$  is

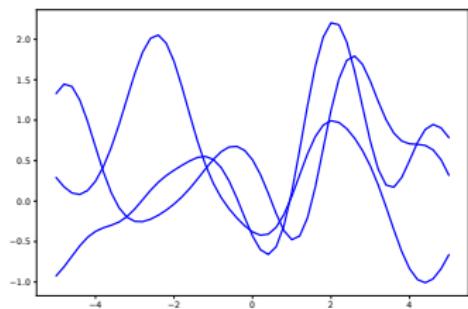
$$\mathbf{K}_{x_*, \mathbf{x}} = [\kappa(x_*, x_1), \dots, \kappa(x_*, x_n)]$$

With Gaussian identity

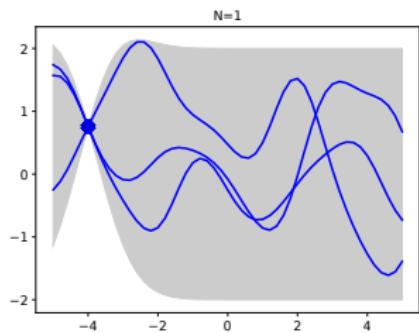
$$\begin{aligned} p(y_* \mid \mathbf{x}_*, \mathcal{D}) &= \mathcal{N}(y_* \mid \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \\ \boldsymbol{\mu}_* &= \mathbf{0} + \mathbf{K}_{\mathbf{x}, x_*}^\top \mathbf{K}_{\mathbf{x}, \mathbf{x}}^{-1} (\mathbf{y} + \mathbf{0}) \\ \boldsymbol{\Sigma}_* &= \mathbf{K}_{x_*, x_*} - \mathbf{K}_{x_*, \mathbf{x}}^\top \mathbf{K}_{\mathbf{x}, \mathbf{x}}^{-1} \mathbf{K}_{\mathbf{x}, x_*} \end{aligned}$$

We can generalize this to multiple predicted values  $\mathbf{y}_*$  given  $\mathbf{x}_*$

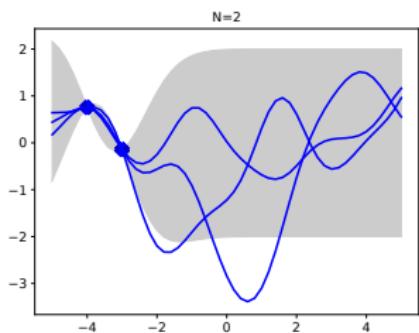
## Examples: Prediction with Observations



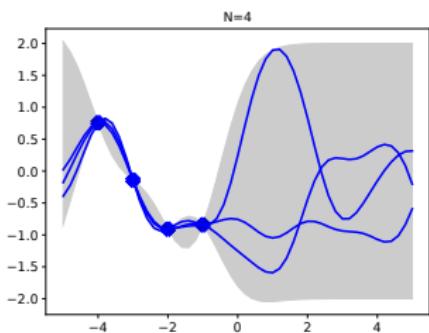
(a)



(b)



(c)



(d)

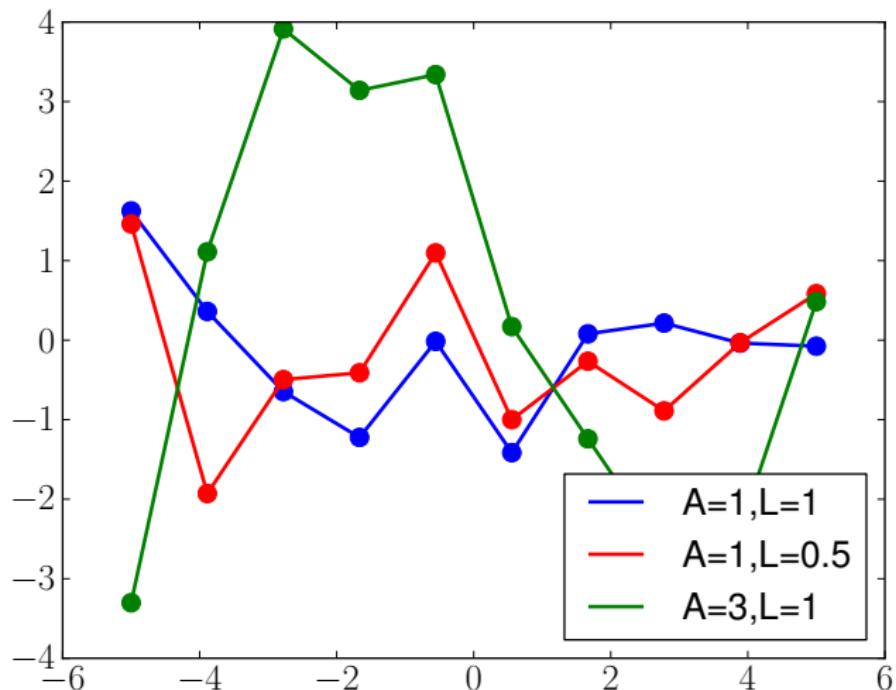
## What does it mean by training Gaussian Process?

A commonly used kernel function to construct covariance is the squared exponential:

$$\kappa_{\text{SE}}(x_i, x_j; A, L) = A^2 \exp \left\{ -0.5 \cdot \frac{(x_i - x_j)^2}{L^2} \right\}$$

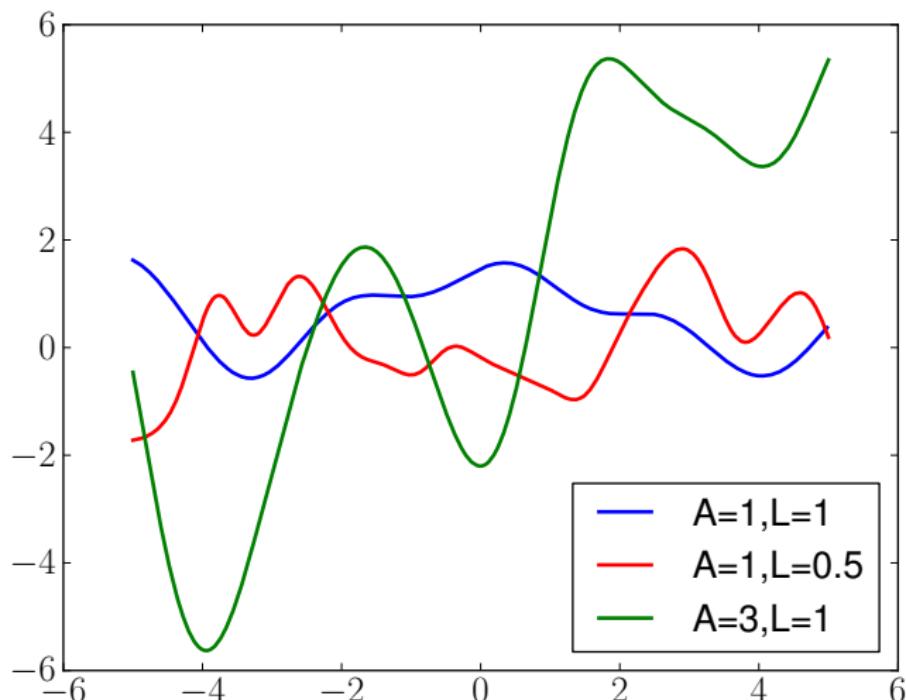
- ▶  $A$ : Overall correlation, amplitude
- ▶  $L$ : Scaling parameter, smoothness
- ▶ Hyperparameters can be tuned to maximize some objective (MLE, MAP)
- ▶ GP aims to find the best kernel function and their hyper-parameters

## Examples: Sampled Functions



10D Gaussian

## Examples: Sampled Functions



500D Gaussian

## Noisy Gaussian Process Regression

- ▶ So far, we have assumed that the observed outputs  $\mathbf{y}$  do not contain noise.
- ▶ As in linear regression, it is more natural to add observation noise into GP:

$$y = f(x) + \epsilon$$

where  $f \sim \text{GP}(m(x), k(x, x'))$  and  $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$ .

- ▶ GP: prior of  $f$
- ▶  $\mathcal{N}(f(x), \sigma_\epsilon^2)$ : likelihood of  $y$
- ▶ We can compute posterior exactly! - Check the textbook

## Further Readings

- ▶ Textbook: Chapter 17.1 (Mercer kernels)
- ▶ Textbook: Chapter 17.2 (Gaussian Processes)
- ▶ Properties of kernels (on PLMS)