



國立交通大學
National Chiao Tung University

Kernel, Gaussian Process, SVM ML 2021 Fall

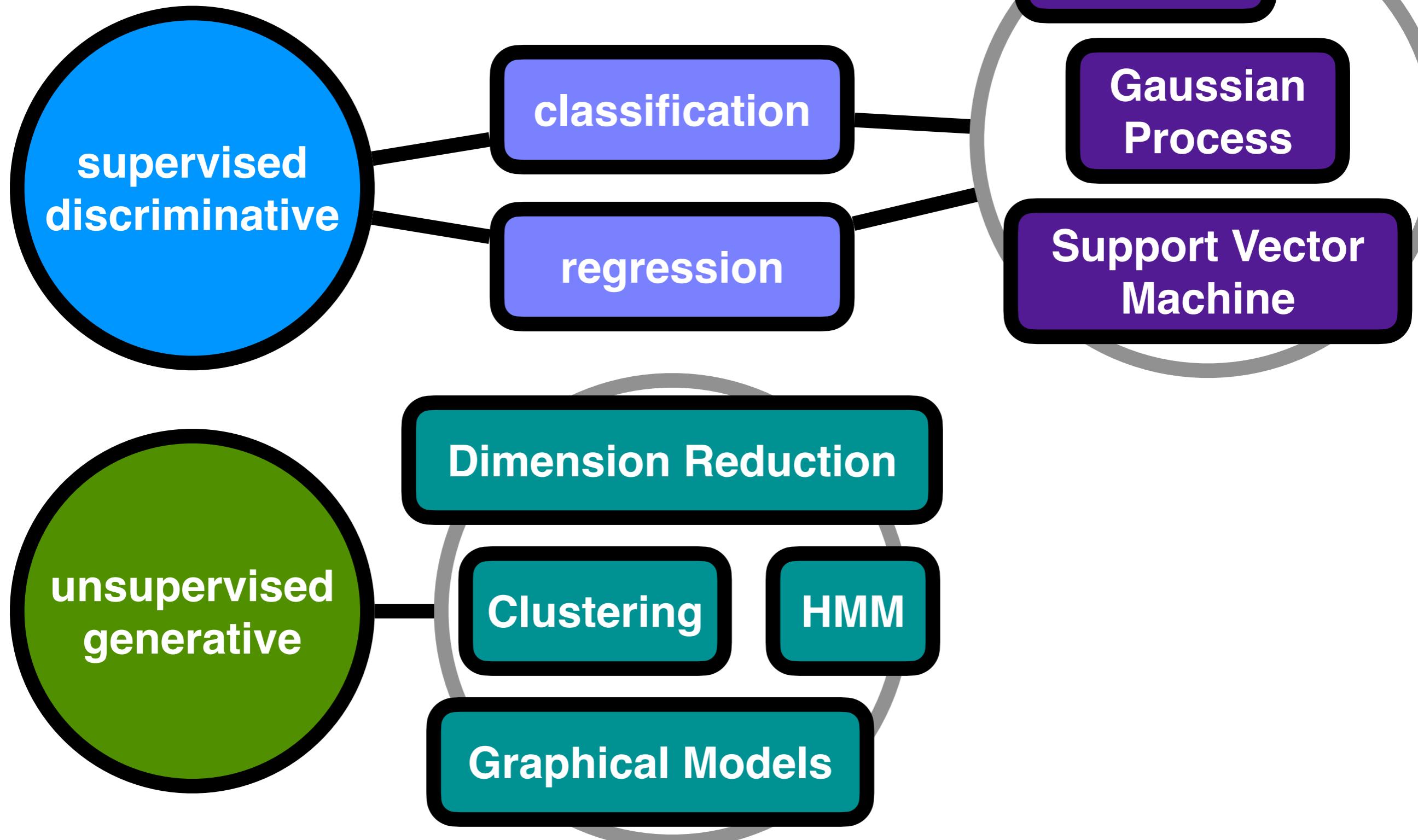
Wei-Chen (Walton) Chiu
邱維辰

Enriched Vision Applications
Laboratory

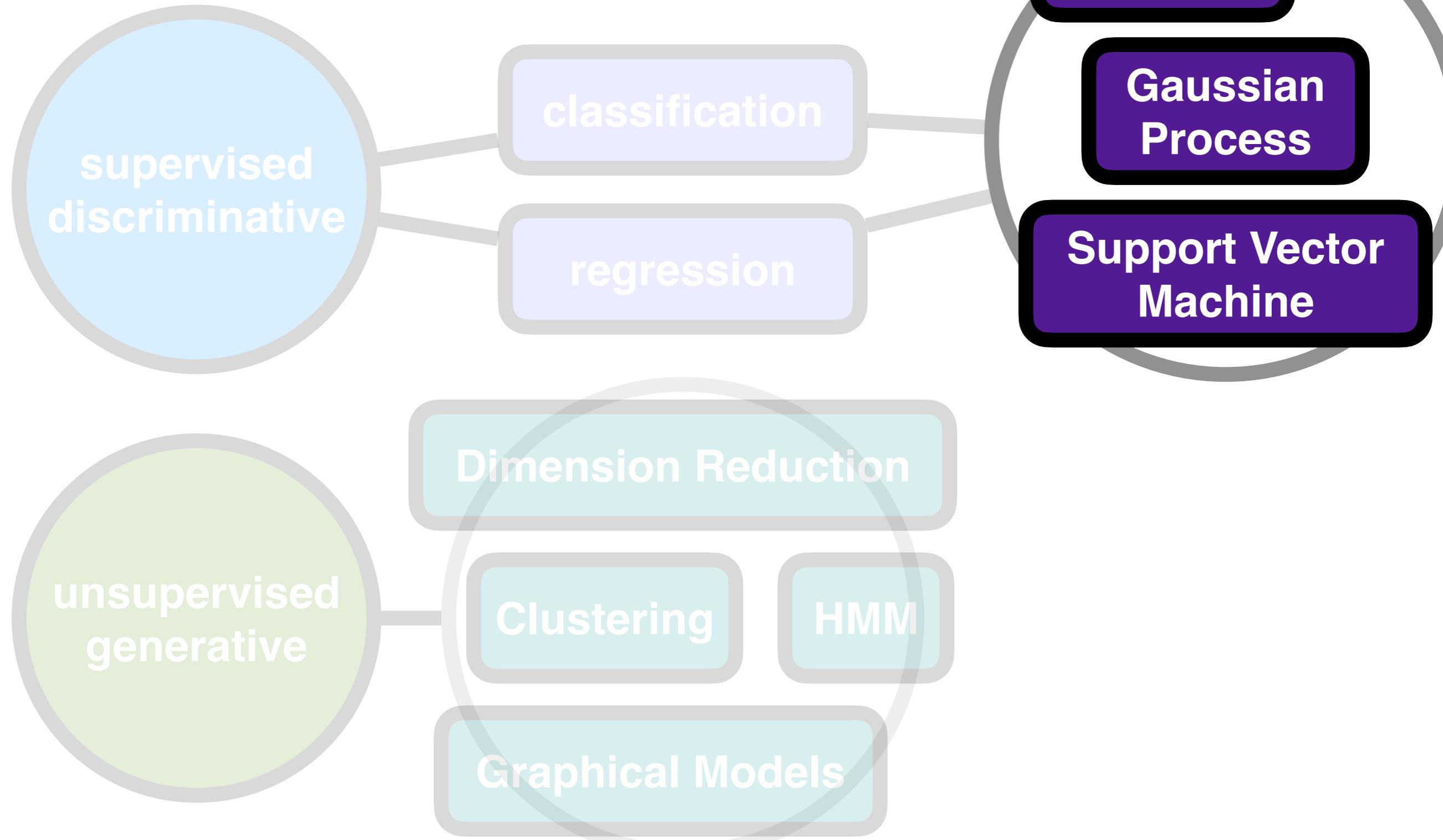


Topics to be tackled around the corner

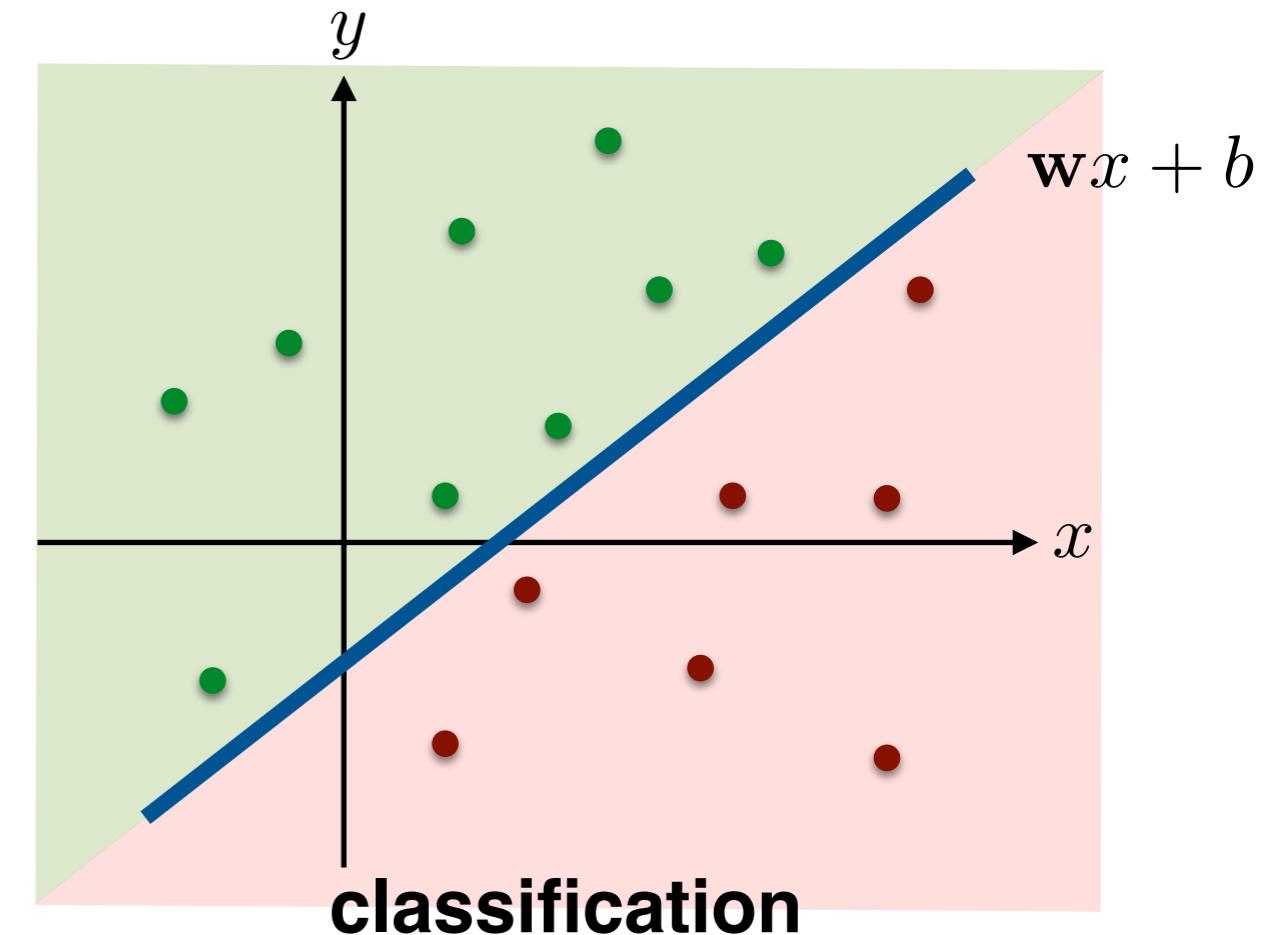
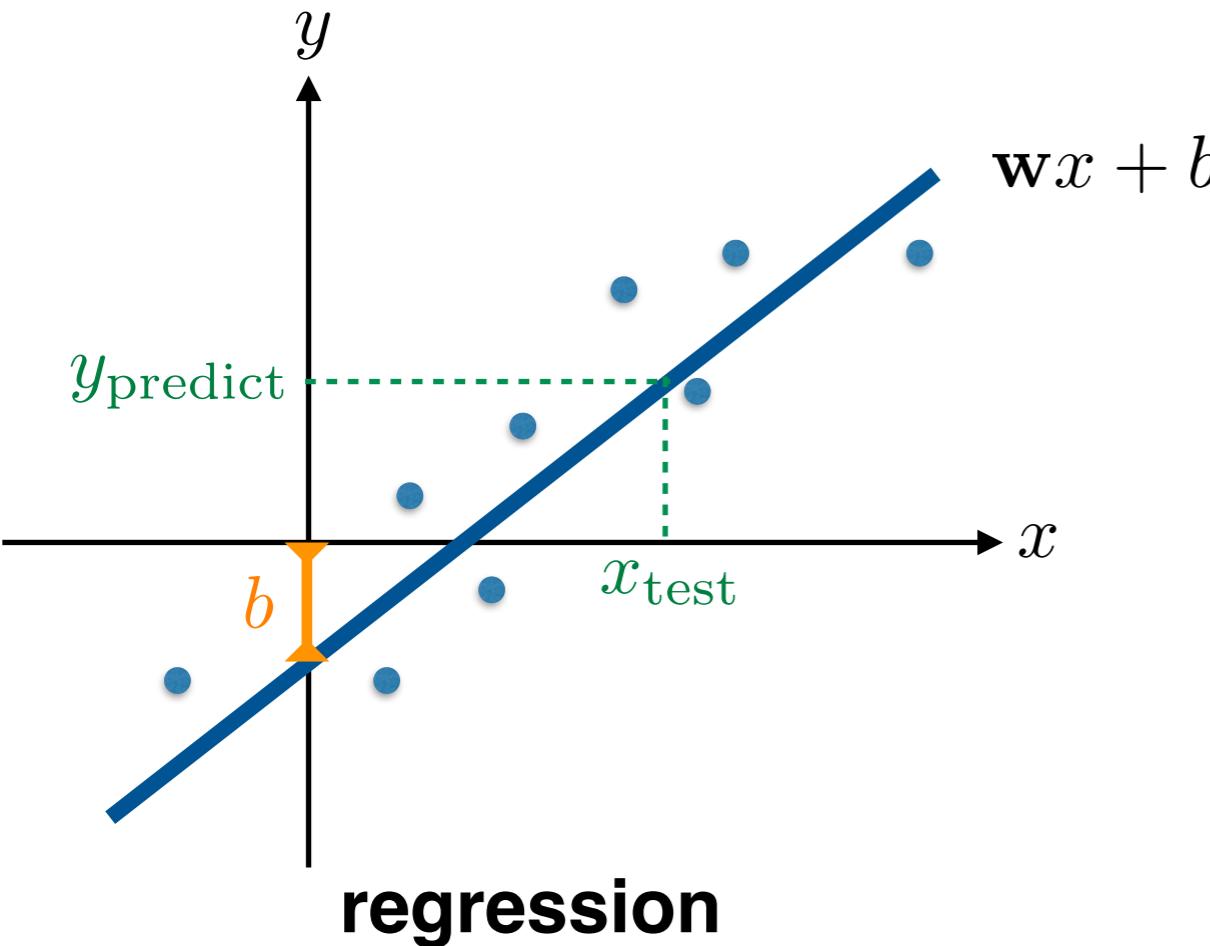
memory-based methods



Today



What you guys have learnt



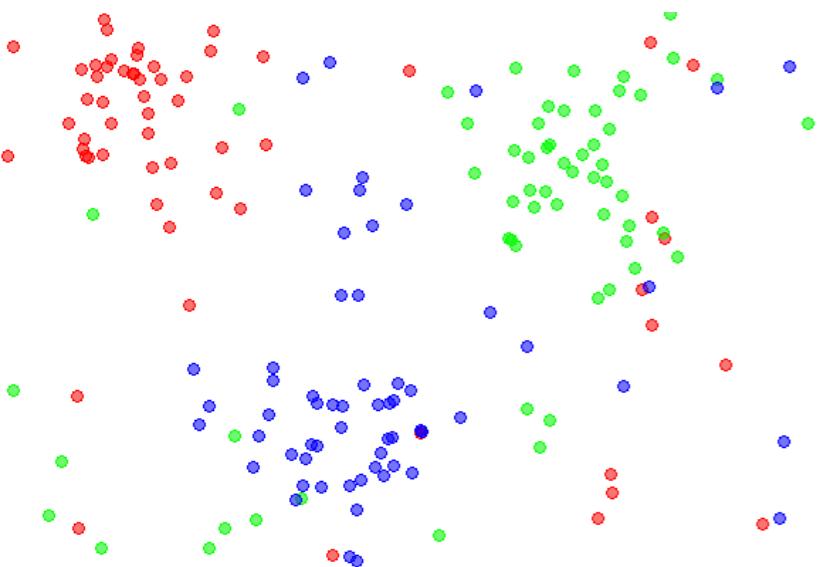
- No matter linear parametric models for regression or classification, throw away the training data after learning the parameters of the model.
- Make prediction based on those parameters only.

Another way of thinking

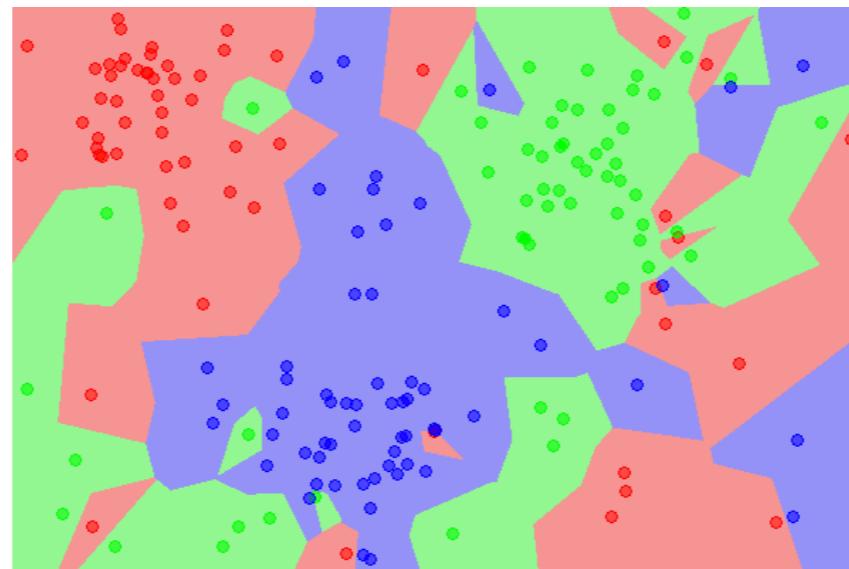
- Training data is valuable, try to keep it even till prediction
 - **Memory-based methods**
 - **We would like to utilise the training data for prediction!**
- For instance, KNN, Parzen

Another way of thinking

- Training data is valuable, try to keep it even till prediction
 - **Memory-based methods**
 - **We would like to utilise the training data for prediction!**
- For instance, **KNN**, Parzen probability estimation



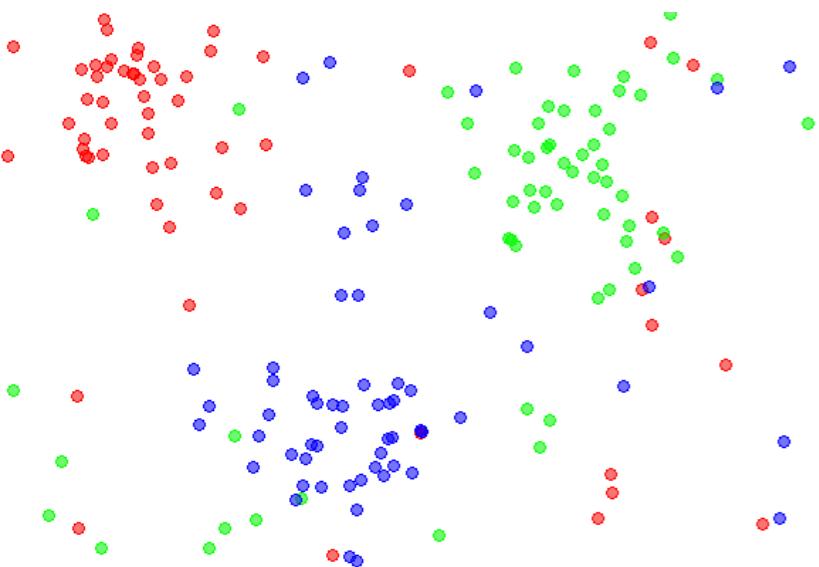
dataset



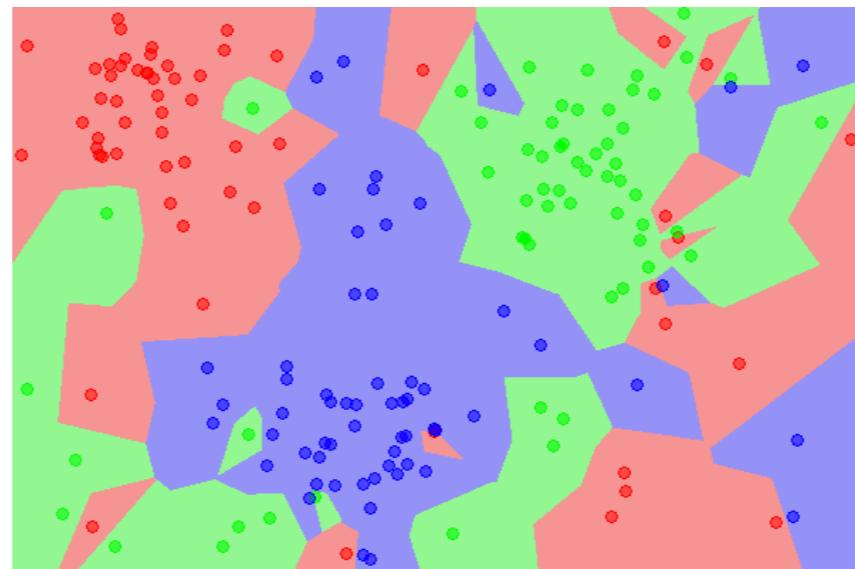
NN (nearest neighbour):
find the closest training data point,
take its class as prediction

Another way of thinking

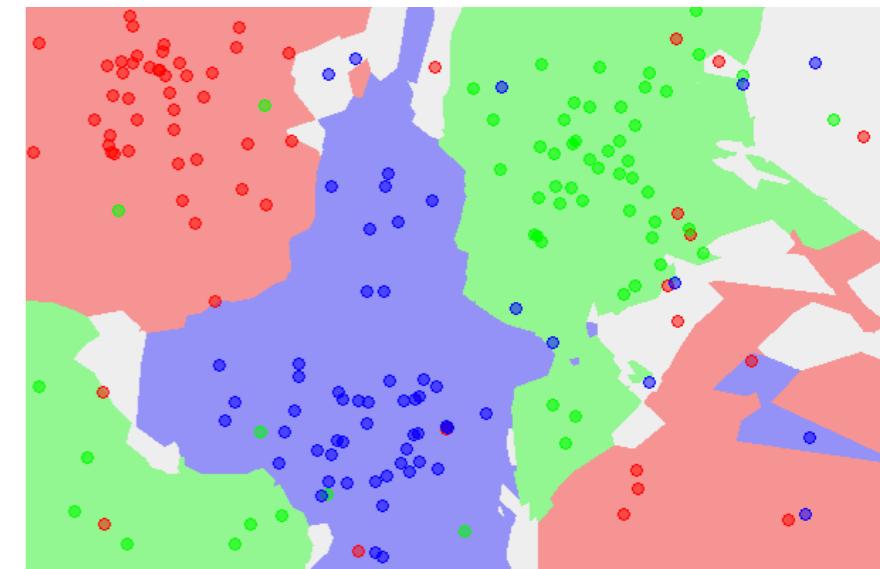
- Training data is valuable, try to keep it even till prediction
 - **Memory-based methods**
 - **We would like to utilise the training data for prediction!**
- For instance, **KNN**, Parzen probability estimation



dataset



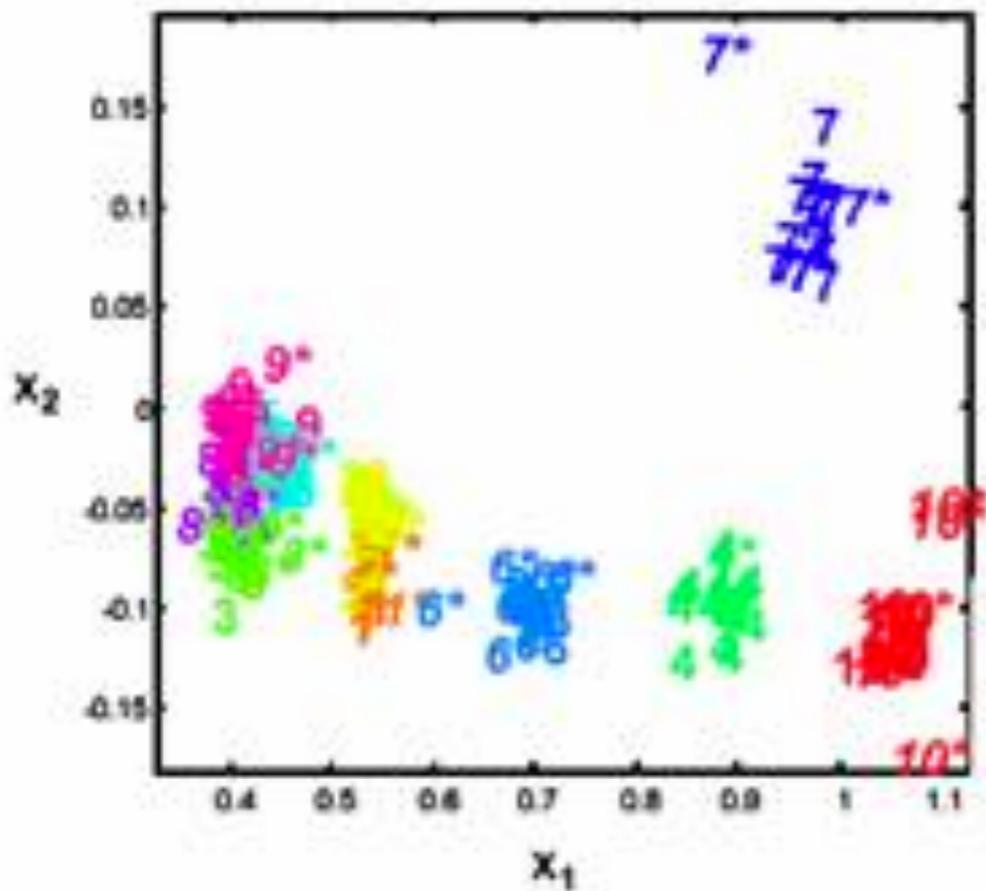
NN (nearest neighbour):



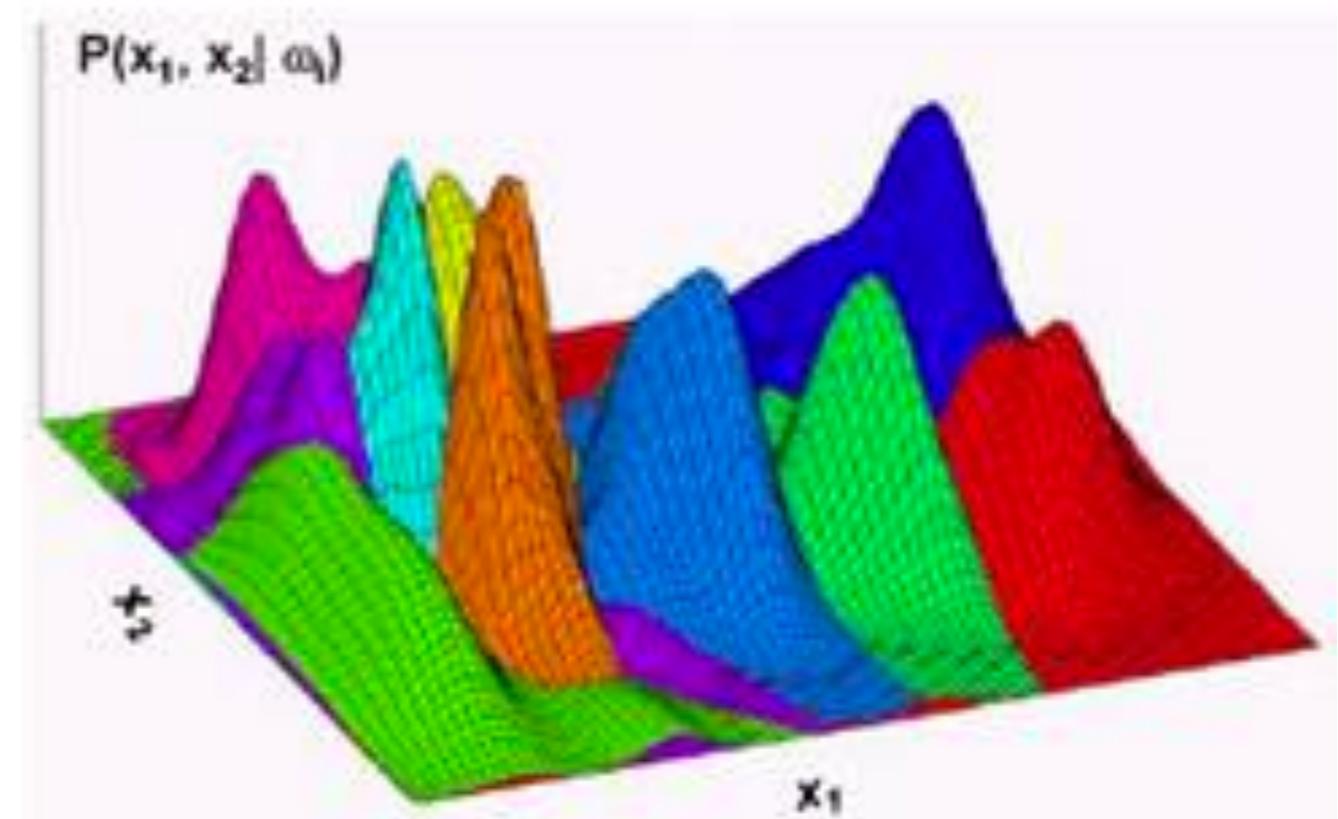
KNN:
prediction based on
K closest training data
points, do majority vote

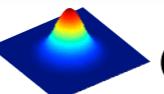
Another way of thinking

- Training data is valuable, try to keep it even till prediction
 - **Memory-based methods**
 - **We would like to utilise the training data for prediction!**
- For instance, KNN, **Parzen probability estimation**



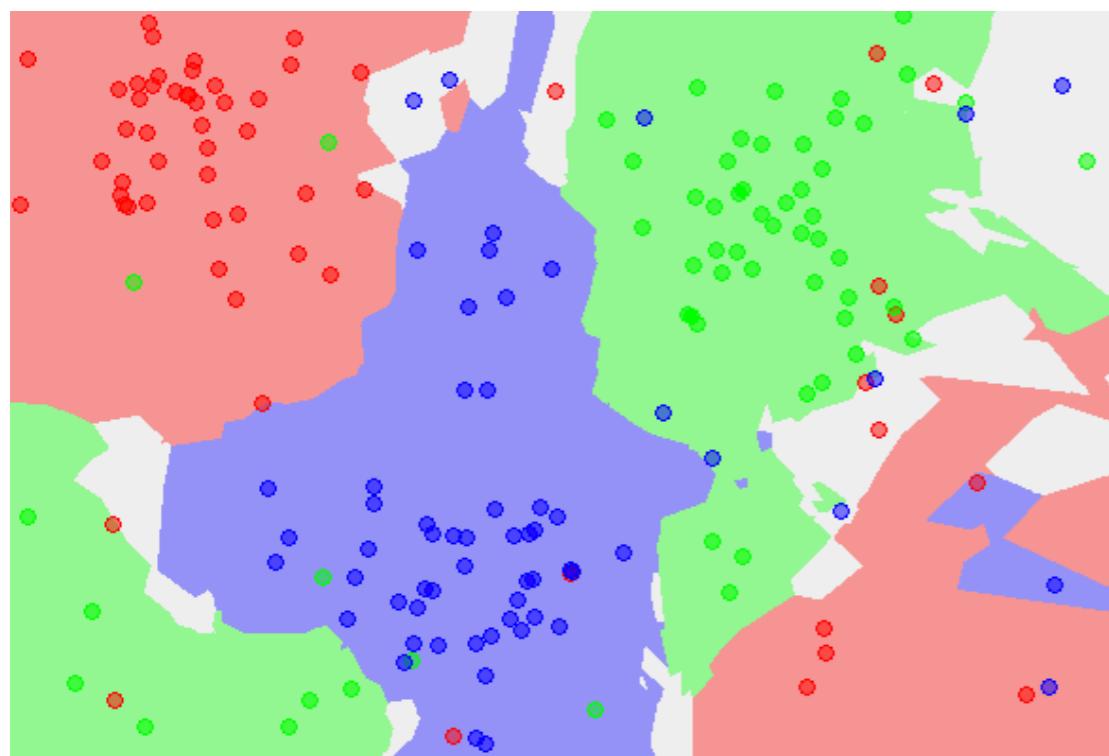
dataset



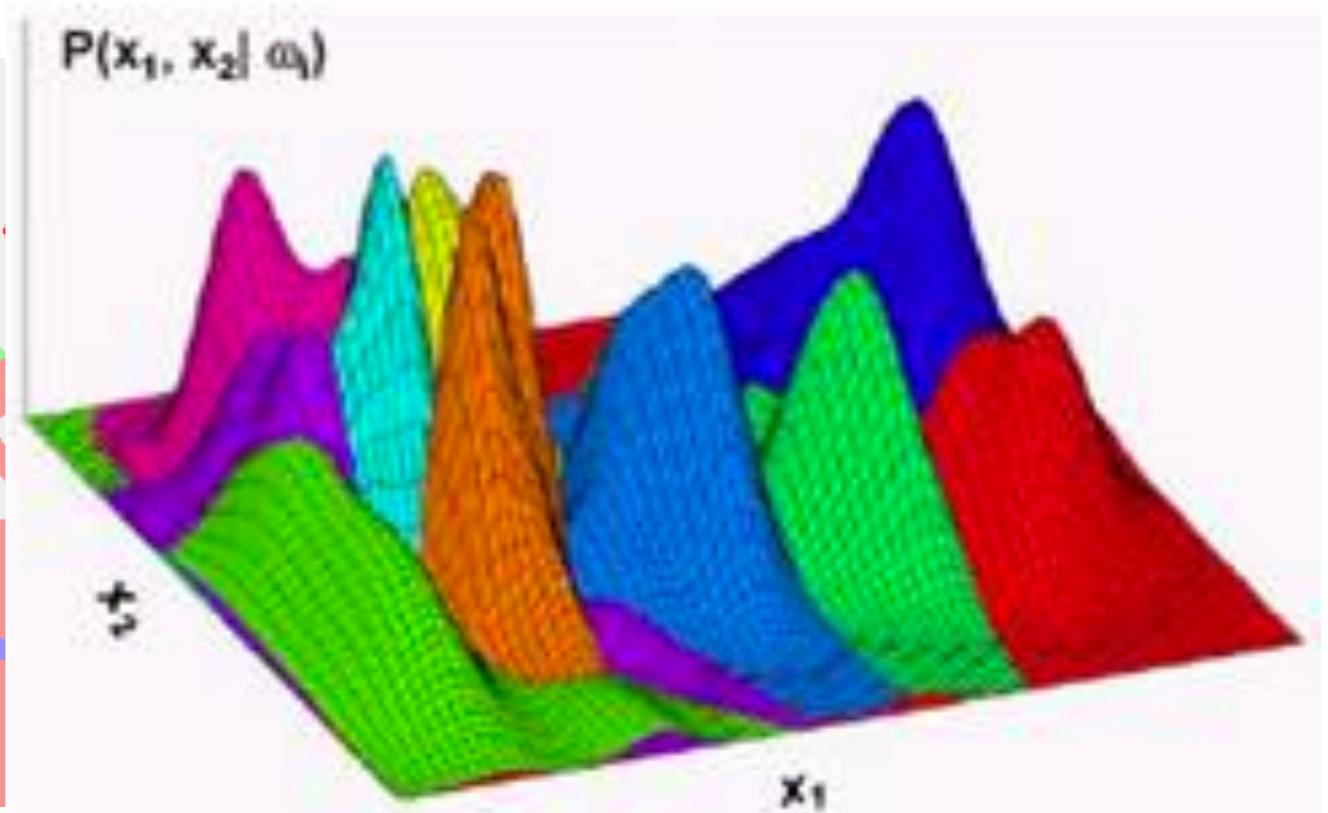
put  on each training data points

Another way of thinking

- Training data is valuable, try to keep it even till prediction
 - **Memory-based methods**
 - **We would like to utilise the training data for prediction!**
 - Need to evaluate the **distance from “testing” to “training” data**



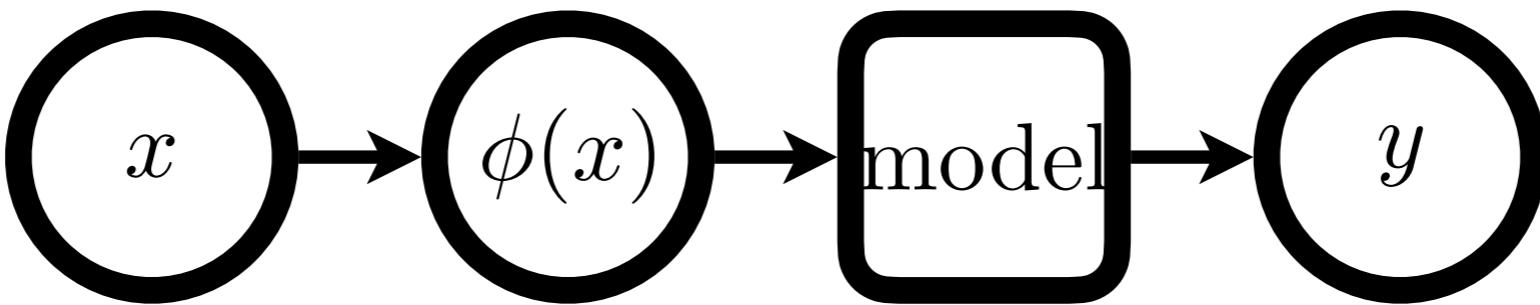
KNN



Parzen probability estimation

Data in the feature space

- General machine learning scheme



- project x to a feature space by feature mapping ϕ
- we need to evaluate the distance/similarity between data $\phi(x)$ in the feature space **dot product** is somehow related to similarity

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

- kernel** function!

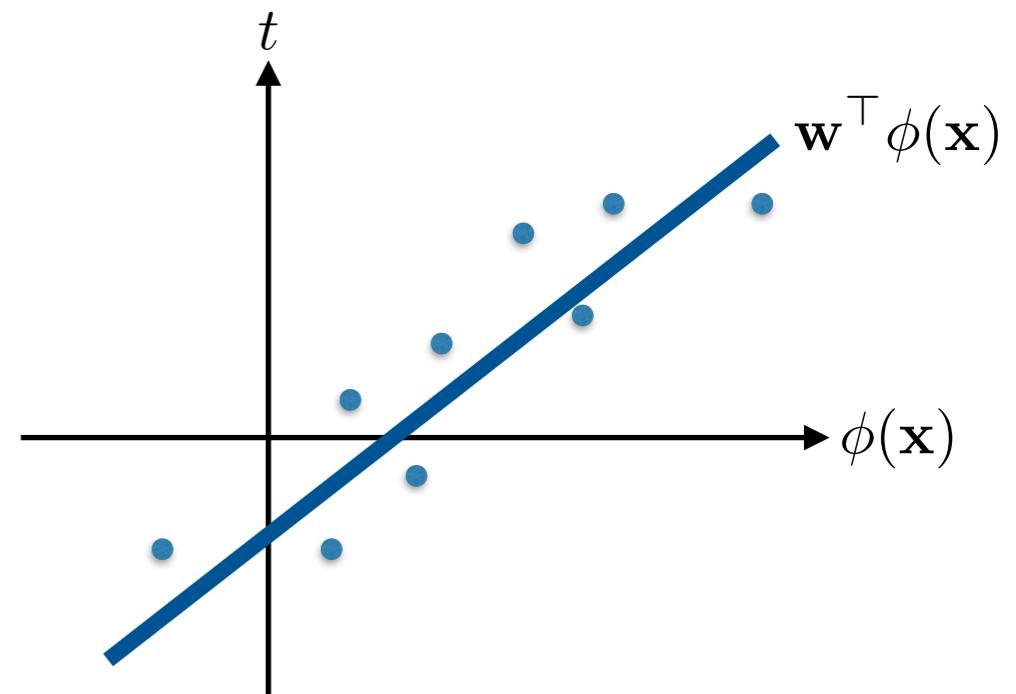
$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$$

- ... so what? why do we need this?

Example for kernel to exist

- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^\top \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$



$$t$$

$$\mathbf{w}^\top \phi(\mathbf{x})$$

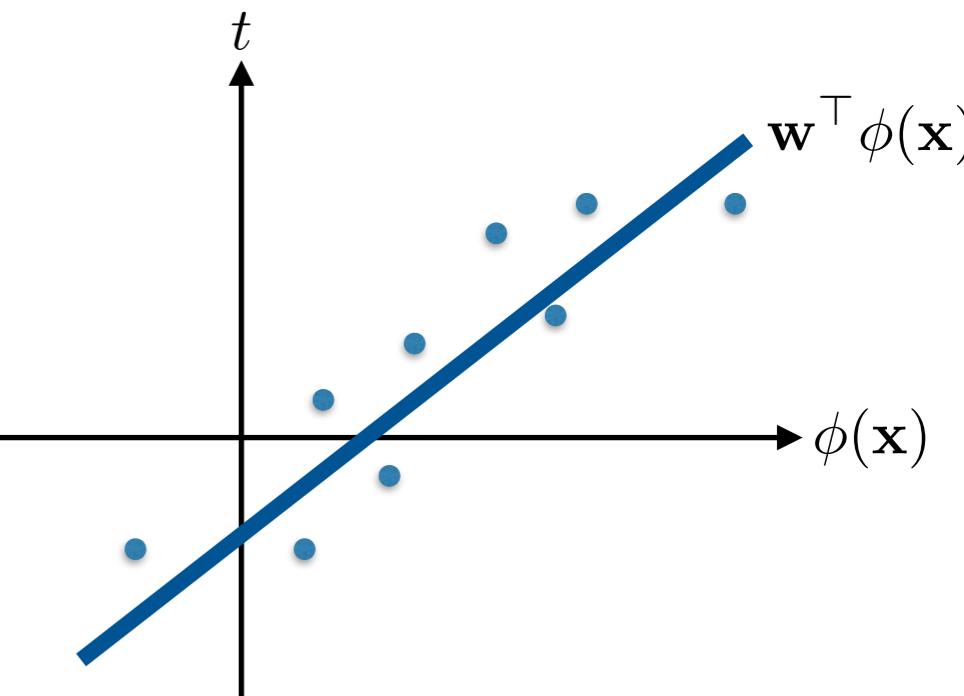
Example for kernel to exist

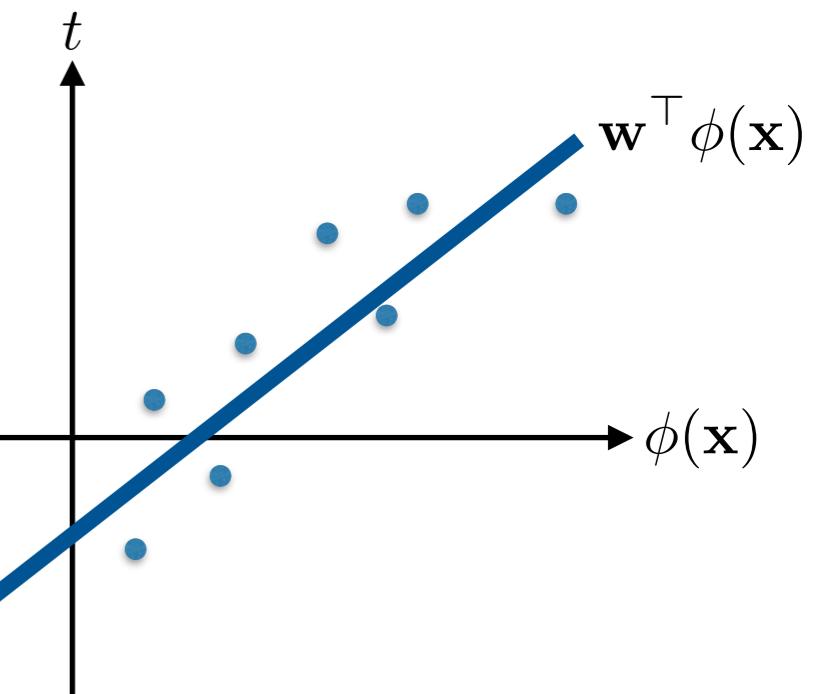
- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^\top \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- with $\frac{\partial J}{\partial \mathbf{w}} = 0$, we can see that \mathbf{w} is the linear combination of $\phi(\mathbf{x})$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \underbrace{\{\mathbf{w}^\top \phi(\mathbf{x}_n) - t_n\}}_{\text{coefficients}} \phi(\mathbf{x}_n) = \sum_{n=1}^N \underbrace{a_n}_{\text{---}} \phi(\mathbf{x}_n) = \Phi^\top \mathbf{a}$$





Example for kernel to exist

- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

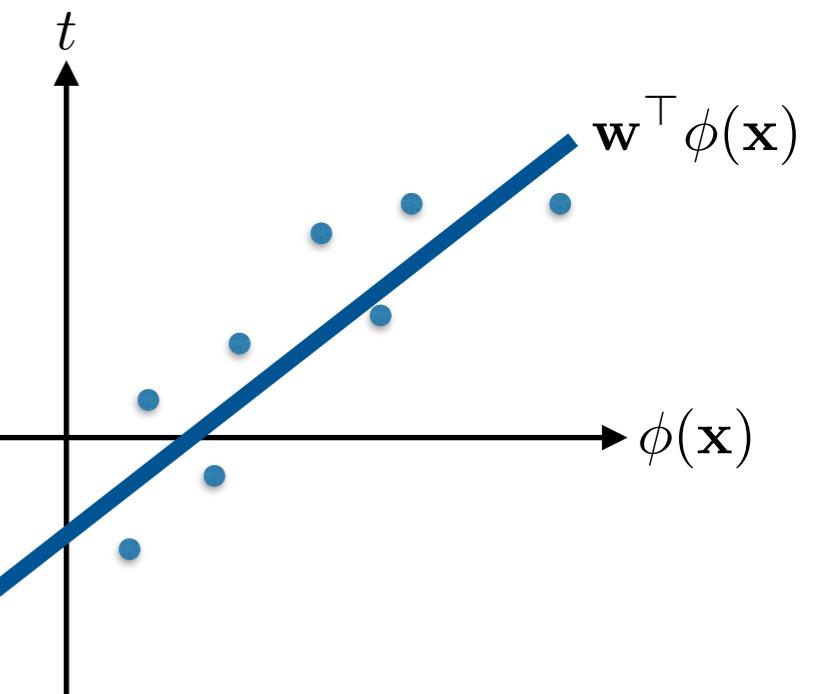
- with $\frac{\partial J}{\partial \mathbf{w}} = 0$, we can see that \mathbf{w} is the linear combination of $\phi(\mathbf{x})$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \underbrace{\{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}}_{\text{coefficients}} \phi(\mathbf{x}_n) = \sum_{n=1}^N \underbrace{a_n}_{\text{---}} \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

- put $\mathbf{w} = \Phi^T \mathbf{a}$ back to J

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \underbrace{\Phi \Phi^T \Phi \Phi^T}_{\text{---}} \mathbf{a} - \mathbf{a}^T \underbrace{\Phi \Phi^T}_{\text{---}} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \underbrace{\Phi \Phi^T}_{\text{---}} \mathbf{a}$$

helloworld! kernel $\phi(\mathbf{x})^\top \phi(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$



Example for kernel to exist

- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$



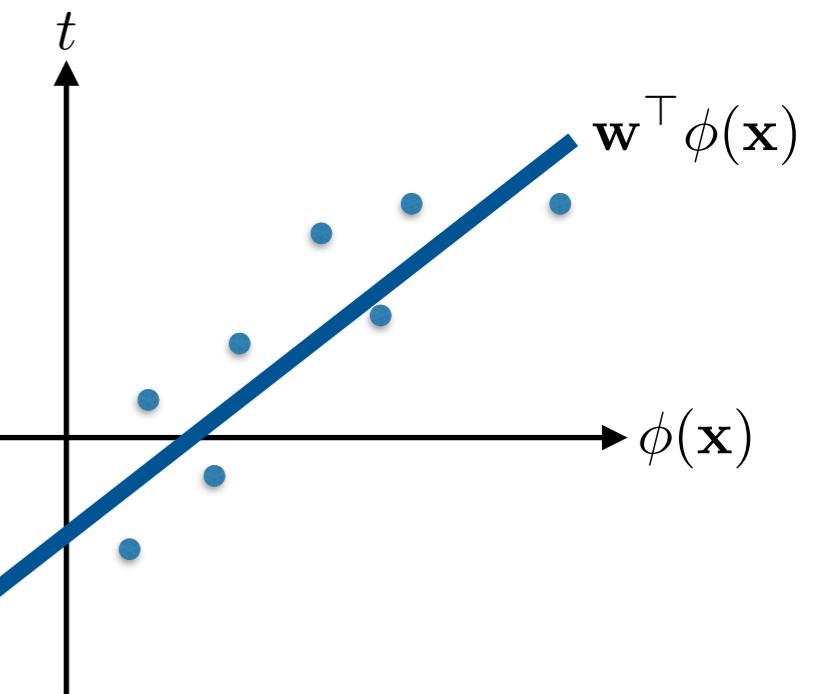
here comes dual form!

- put $\mathbf{w} = \Phi^\top \mathbf{a}$ back to J

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \underline{\Phi \Phi^\top \Phi \Phi^\top} \mathbf{a} - \mathbf{a}^T \underline{\Phi \Phi^\top} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \underline{\Phi \Phi^\top} \mathbf{a}$$

- $K_{nm} = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$ Gram matrix: $K = \Phi \Phi^\top$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T K K \mathbf{a} - \mathbf{a}^T K \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T K \mathbf{a}.$$



Example for kernel to exist

- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

dual form

- $J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}.$
- with $\frac{\partial J}{\partial \mathbf{a}} = 0$, get $\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$.
- when performing prediction:

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

where vector $\mathbf{k}(\mathbf{x})$ with elements $k_n(\mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x})$

$$t$$

$$\mathbf{w}^\top \phi(\mathbf{x})$$

Example for kernel to exist

- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^\top \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

dual form

- $$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^\top \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^\top \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^\top \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^\top \mathbf{K} \mathbf{a}.$$

- with $\frac{\partial J}{\partial \mathbf{a}} = 0$, get $\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$.

- when performing prediction:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) = \mathbf{a}^\top \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

where vector $\mathbf{k}(\mathbf{x})$ with elements $k_n(\mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x})$

**prediction based on
linear combination of
kernel functions evaluated
at training data points**



Example for kernel to exist

- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

dual form

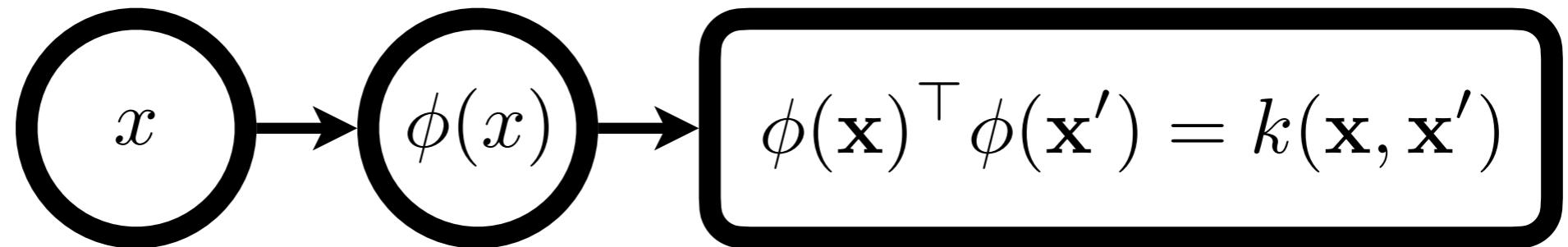
$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}.$$



- with $\frac{\partial J}{\partial \mathbf{a}} = 0$, get $\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$. **if x always appears as scalar product, it might be good to try kernel! known as “kernel trick” a.k.a. “kernel substitution”**

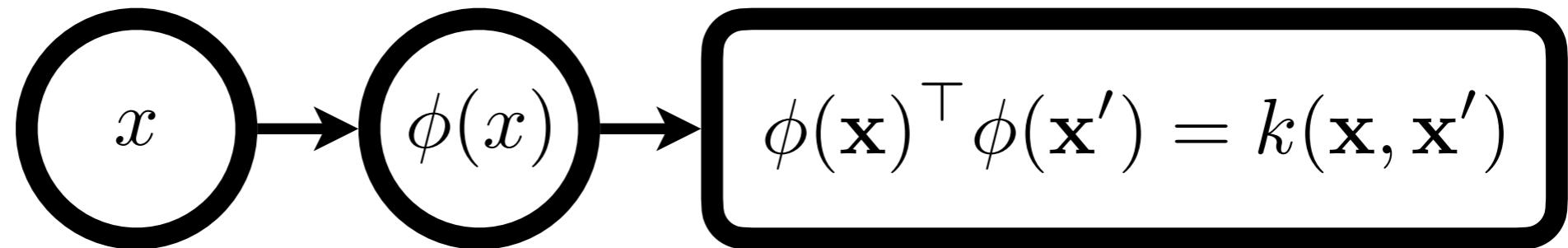
Feature, feature

- All about features:
 - **projection from original data space to feature space**
 - How to choose a proper feature space?
 - for instance: easier separation between classes (classification)
 - prior knowledge about the class of functions to be learned
 - subsets of a basis of the function space (e.g. Fourier, Wavelet, etc.)

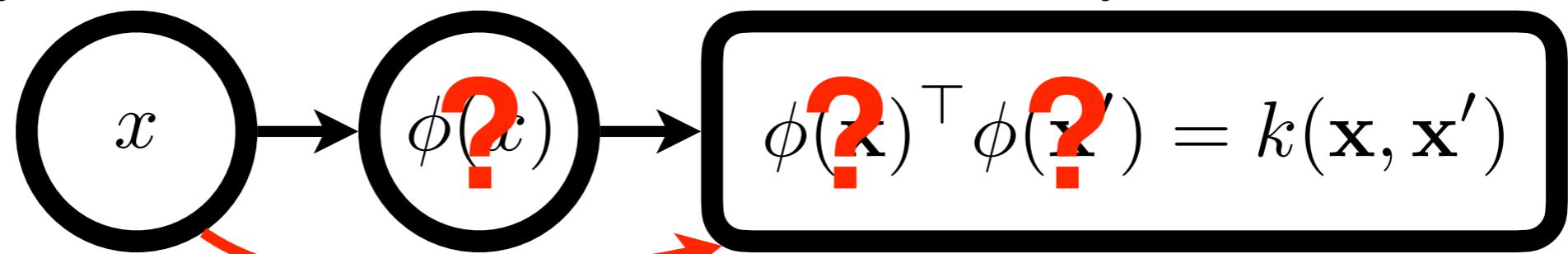


Feature, feature

- All about features:
 - **projection from original data space to feature space**
 - How to choose a proper feature space?
 - for instance: easier separation between classes (classification)
 - prior knowledge about the class of functions to be learned
 - subsets of a basis of the function space (e.g. Fourier, Wavelet, etc.)

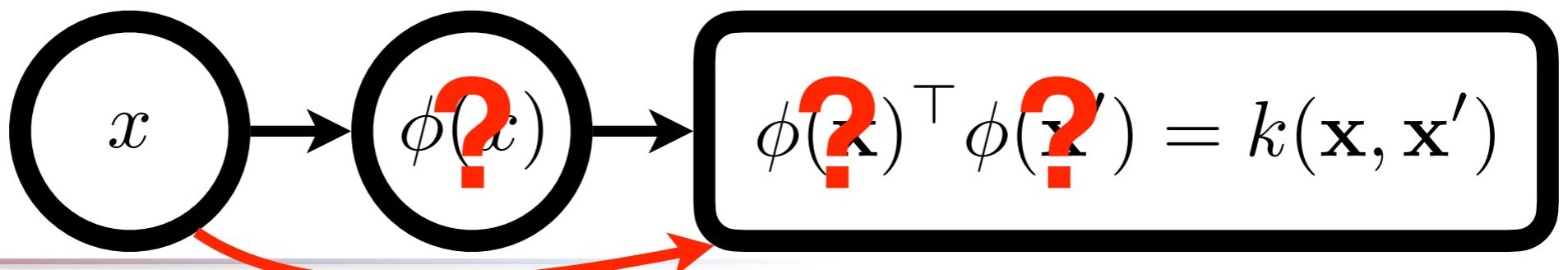


- Sometimes we don't know what is the proper feature mapping...
 - still want to use memory-based methods (training data in prediction)
 - perhaps just some kernel? don't care what exactly is the feature map?



Some examples of kernel functions

- a very simple form $k(x, x') = e^{xx'}$



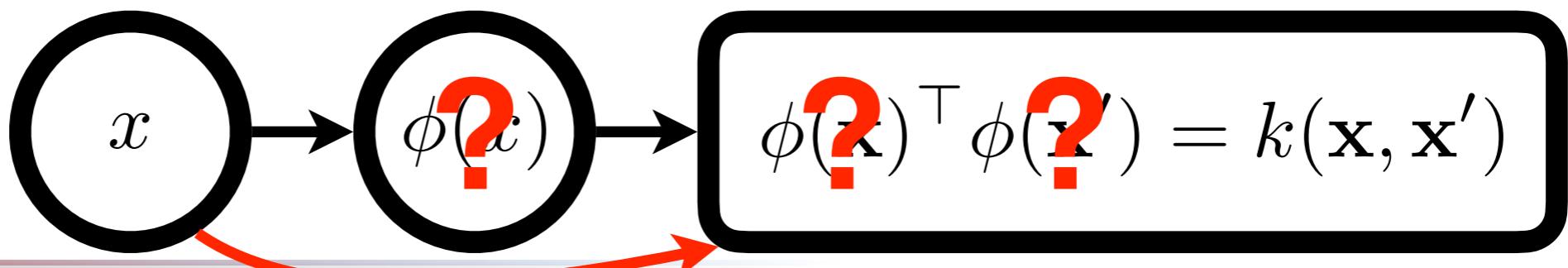
Some examples of kernel functions

- a very simple form $k(x, x') = e^{xx'}$
 - ▶ the corresponding feature map is

$$\forall x \in \mathbb{R}, r \in \mathbb{N} : \phi_r(x) = \frac{1}{\sqrt{r!}} x^r$$

$$k(x, x') = \sum_{r=0}^{\infty} \phi_r(x) \phi_r(x') = \sum_{r=0}^{\infty} \frac{x^r}{\sqrt{r!}} \frac{(x')^r}{\sqrt{r!}} = \sum_{r=0}^{\infty} \frac{(xx')^r}{r!} = e^{xx'}$$

***simple kernel but it has even
countably infinitely many feature maps!***



Some examples of kernel functions

- a very simple form $k(x, x') = e^{xx'}$
 - ▶ the corresponding feature map is

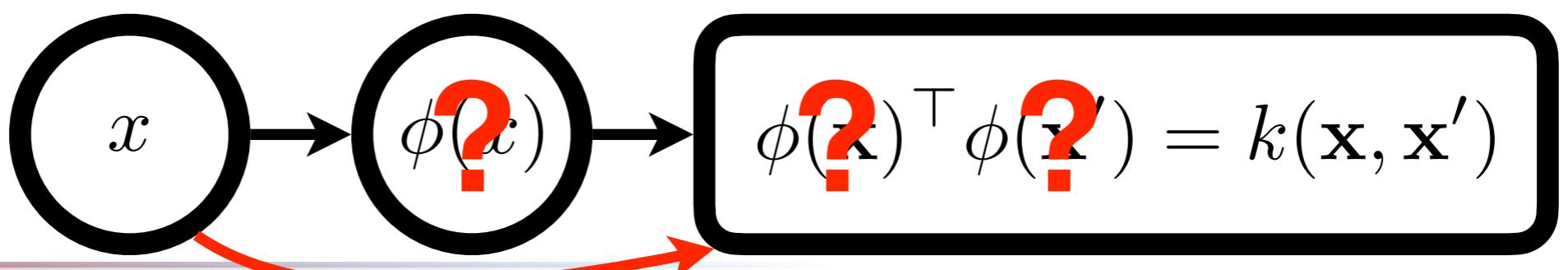
$$\forall x \in \mathbb{R}, r \in \mathbb{N} : \phi_r(x) = \frac{1}{\sqrt{r!}} x^r$$

$$k(x, x') = \sum_{r=0}^{\infty} \phi_r(x) \phi_r(x') = \sum_{r=0}^{\infty} \frac{x^r}{\sqrt{r!}} \frac{(x')^r}{\sqrt{r!}} = \sum_{r=0}^{\infty} \frac{(xx')^r}{r!} = e^{xx'}$$

**simple kernel but it has even
countably infinitely many feature maps!**



PS: simplest kernel is identity mapping
 $\phi(x) = x, k(x, x') = x^\top x'$
it is called linear kernel.



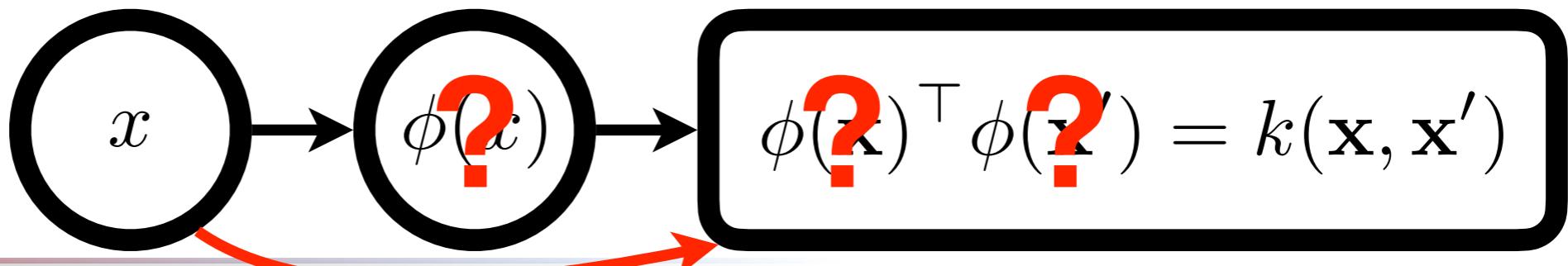
Some examples of kernel functions

- Given $x_i \in \mathbb{R}^3, \phi(x_i) \in \mathbb{R}^{10}$

$$\phi(x_i) = \left[1, \sqrt{2}(x_i)_1, \sqrt{2}(x_i)_2, \sqrt{2}(x_i)_3, (x_i)_1^2, (x_i)_2^2, (x_i)_3^2, \right. \\ \left. \sqrt{2}(x_i)_1(x_i)_2, \sqrt{2}(x_i)_1(x_i)_3, \sqrt{2}(x_i)_2(x_i)_3 \right]^\top$$

- kernel function: $\phi(x_i)^\top \phi(x_j) = \underline{(1 + x_i^\top x_j)^2}$

**feature map projects 3-D to 10-D, but it results as
a very simple form to compute kernel!**



Some examples of kernel functions

- Radial Basis Kernel $e^{-\gamma \|x_i - x_j\|}$

► assume $x \in \mathbb{R}^1$ and $\gamma > 0$

$$e^{-\gamma \|x_i - x_j\|} = e^{-\gamma(x_i - x_j)^2} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2}$$

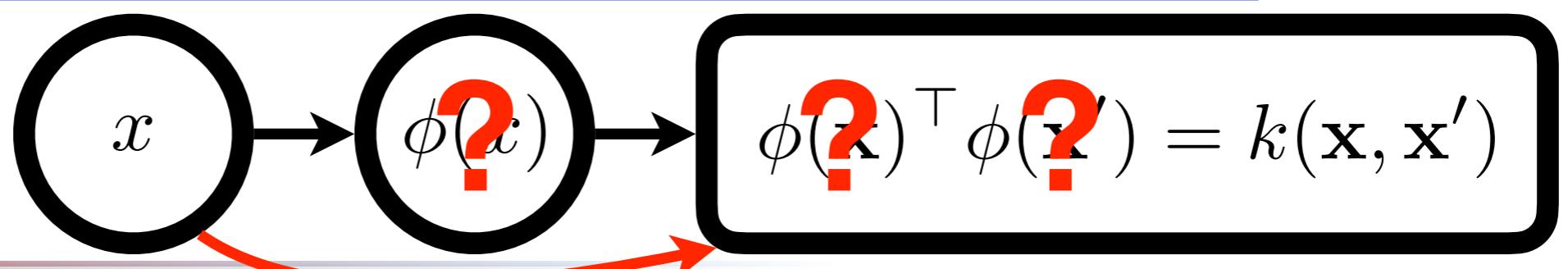
$$= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \dots\right)$$

$$= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}} x_i \times \sqrt{\frac{2\gamma}{1!}} x_j + \sqrt{\frac{(2\gamma)^2}{2!}} x_i^2 \times \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 + \dots\right)$$

$$= \phi(x_i)^\top \phi(x_j)$$

where $\phi(x) = e^{-\gamma x^2} \left[1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots\right]^\top$

simple kernel but infinite-D feature map



Example for kernel to exist

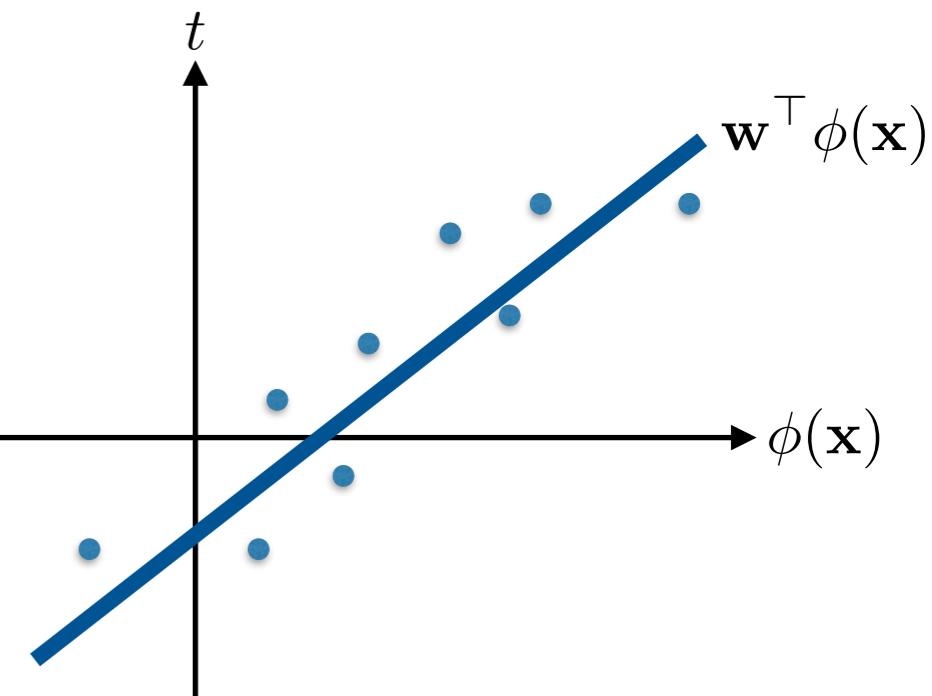
- Linear regression with L_2 regularisation

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

dual form

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}.$$

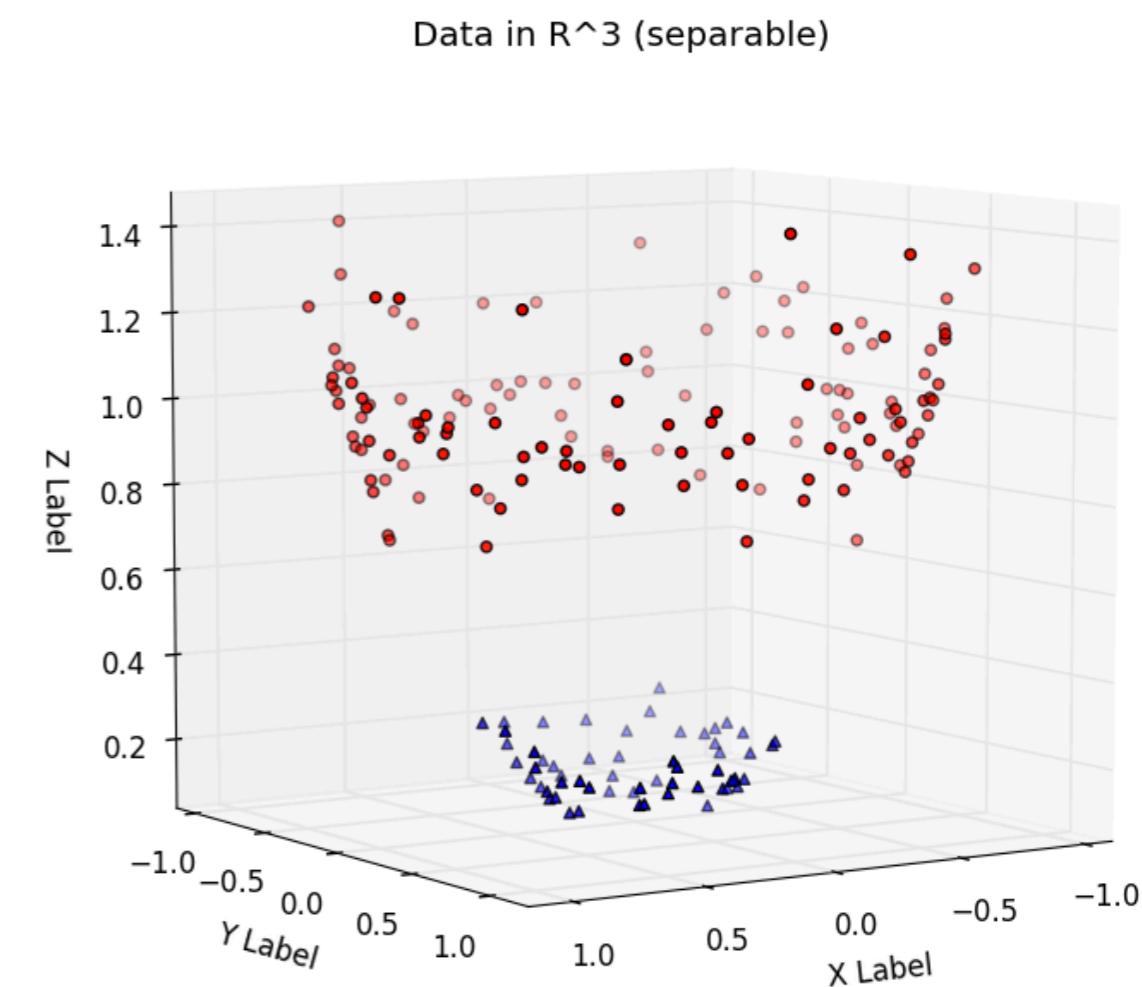
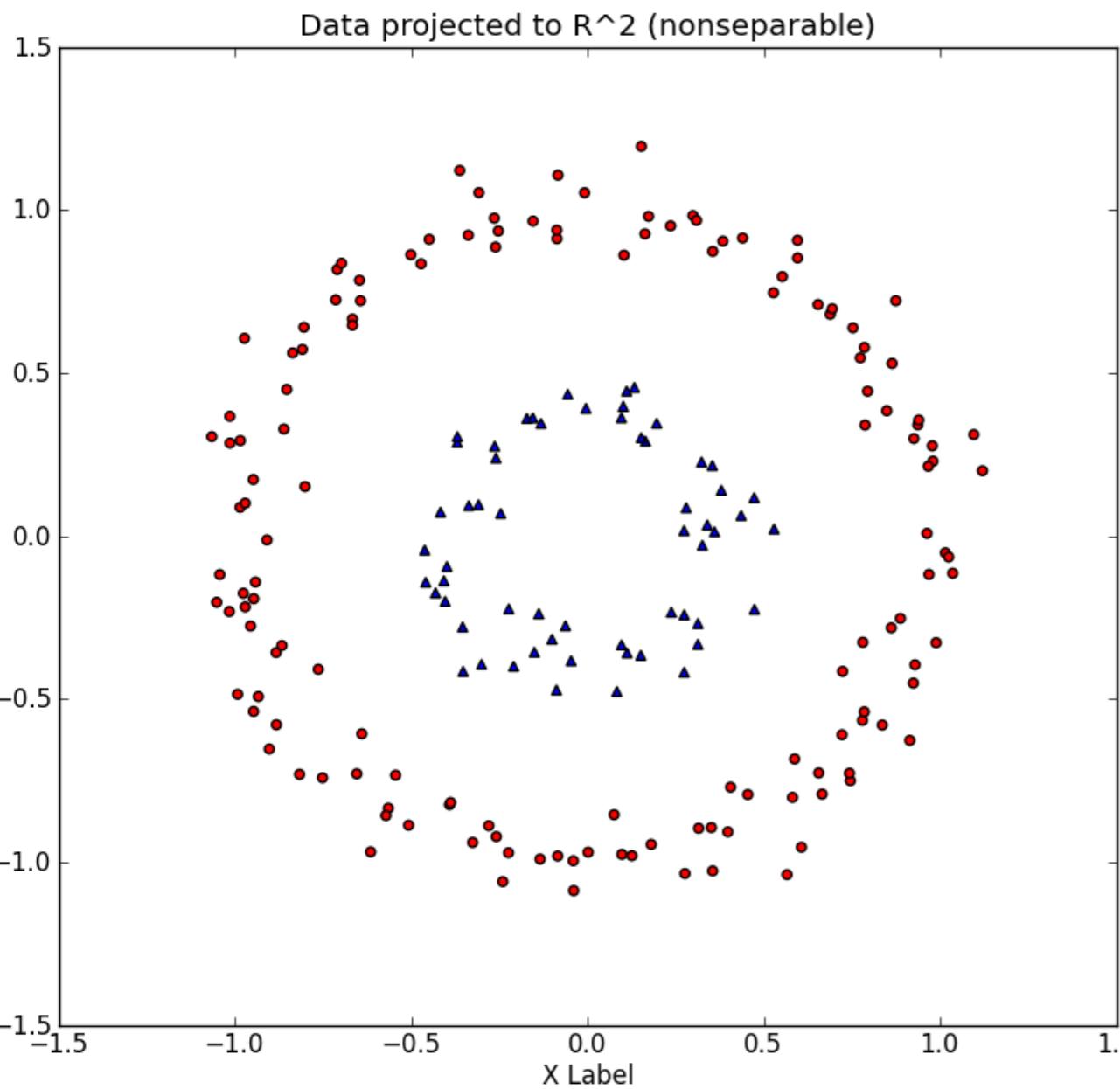
$$\text{with } \frac{\partial J}{\partial \mathbf{a}} = 0, \text{ get } \boxed{\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}.}$$



ah-ha!

even it seems more complicated to do matrix inversion in higher-D space, but we can directly work on kernels (usually simple) and avoid explicit introduction of feature vector $\phi(x)$, which allows us implicitly to use feature spaces of high, even infinite, dimensionality!

Why High Dimensional Feature Space



How to construct valid kernel?

- a necessary and sufficient condition:
the Gram matrix: $K = \Phi\Phi^\top$ whose elements are given by
 $k(x, x') = \phi(x)^\top \phi(x')$ is positive semidefinite

A symmetric function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a **positive definite kernel** if, for any finite J , any $x \in \mathcal{X}^J$ and any $c \in \mathbb{R}^J$:

$$0 \leq \sum_{i \in J} \sum_{j \in J} c_i c_j k(x_i, x_j)$$

The set of all real-valued positive definite kernels on \mathcal{X} is denoted $\mathbb{R}_+^{\mathcal{X} \times \mathcal{X}}$

How to construct valid kernel?

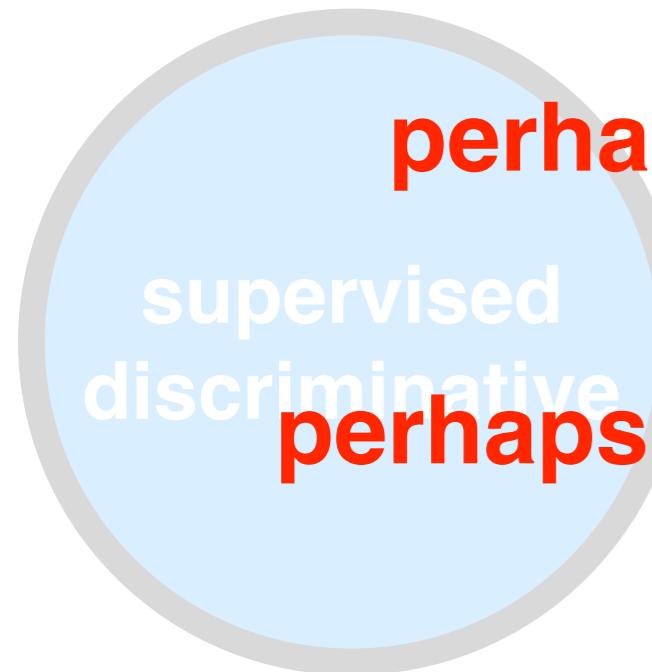
- a necessary and sufficient condition:
the Gram matrix: $K = \Phi\Phi^\top$ whose elements are given by
 $k(x, x') = \phi(x)^\top \phi(x')$ is positive semidefinite

- You can build kernel from kernel
 - $k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}')$
 - $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$
 - $k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$
 - $k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$
 - $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
 - $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$
 - $k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}'))$
 - $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{A} \mathbf{x}'$
 - $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b)$
 - $k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)$

Recap!

- memory-based methods: we want to keep training data points for further use in the prediction
- need to have a metric to evaluate the distance/similarity between data points in the feature space: inner product  kernel!
 - scalar product in \mathbf{x} brings possibility to have kernel trick!
- kernel provide a way to compute inner product in feature space
 - simple kernel computation can be related to complicated feature map!
 - the kernel function itself is more important than the feature map (for memory-based methods), we don't even need to know how to compute feature map!
 - usually we imagine kernel is providing us a way to project the data into much higher dimensional space
 - we can easily build kernel from kernel!

Today



**perhaps most powerful regressor
before deep learning**

**perhaps most powerful classifier
before deep learning**



Dimension Reduction

Clustering

HMM

Graphical Models

Kernel
Methods

Gaussian
Process

Support Vector
Machine

Nonparametric Regression

- Previously on linear parametric regression

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$$

- prior distribution over \mathbf{w} , e.g., isotropic Gaussian $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$, induces a corresponding distribution over functions $y(\mathbf{x}, \mathbf{w})$
- while given training samples with specific values $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, we can get random variables $y(\mathbf{x}_1), y(\mathbf{x}_2), \dots, y(\mathbf{x}_N)$
- the joint distribution $\mathbf{y} = \{y(\mathbf{x}_1), y(\mathbf{x}_2), \dots, y(\mathbf{x}_N)\}$ is also Gaussian!

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^\top] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^\top] \Phi^\top = \frac{1}{\alpha} \Phi \Phi^\top = \mathbf{K}$$

Nonparametric Regression

- Previously on linear parametric regression

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$$

- prior distribution over \mathbf{w} , e.g., isotropic Gaussian $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$, induces a corresponding distribution over functions $y(\mathbf{x}, \mathbf{w})$
- while given training samples with specific values $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, we can get random variables $y(\mathbf{x}_1), y(\mathbf{x}_2), \dots, y(\mathbf{x}_N)$
- the joint distribution $\mathbf{y} = \{y(\mathbf{x}_1), y(\mathbf{x}_2), \dots, y(\mathbf{x}_N)\}$ is also Gaussian!

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^\top] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^\top] \Phi^\top = \frac{1}{\alpha} \Phi \Phi^\top = \mathbf{K}$$

- Now we are going to introduce “**Gaussian Process**”, which doesn’t use parametric model but instead define a prior probability **distribution over functions** directly.

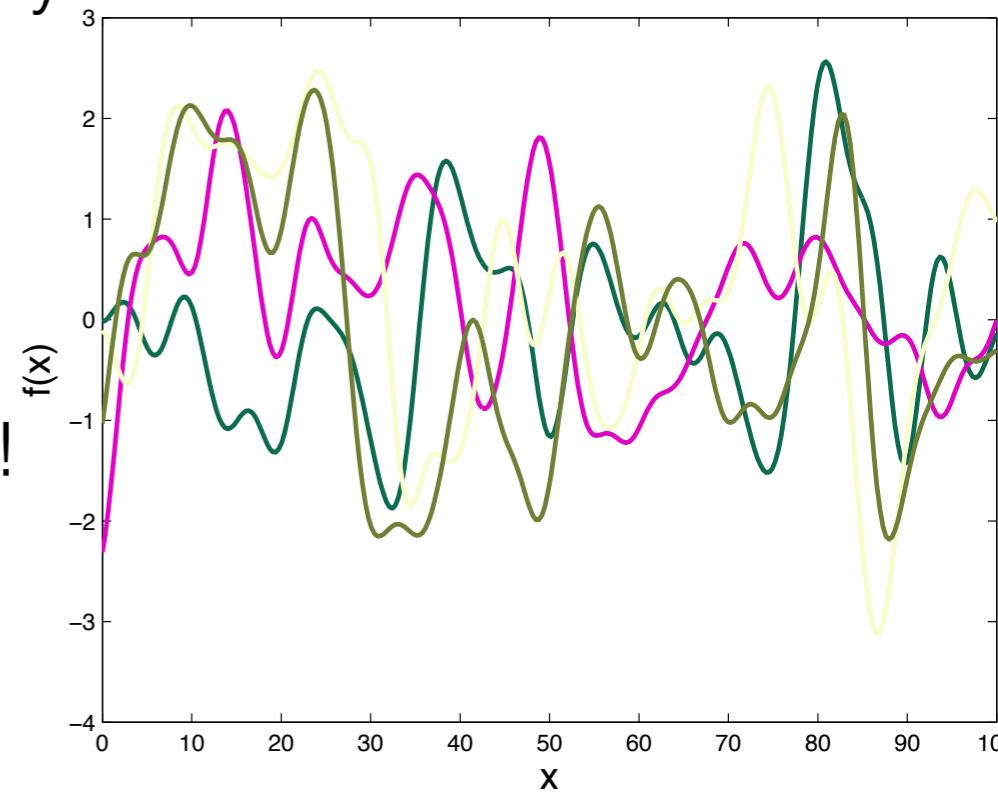


what do these mean?

Gaussian Process

For any set \mathbf{S} , a Gaussian Process on \mathbf{S} is a set of random variables ($f(x), x \in \mathbf{S}$), so for any $n \in \mathbb{N}$ and $x_1, \dots, x_n \in \mathbf{S}, \{f(x_1), \dots, f(x_n)\}$ is (multivariate) Gaussian

- the number of elements in set \mathbf{S} can be infinite many.
- the number of random variables $f(x)$ can be any
 - $\{f(x_9), f(x_5), f(x_2), f(x_7)\}$ can build a mean and a covariance function, Gaussian!
 - $\{f(x_7), f(x_8)\}$ can build another mean and another covariance function, Gaussian!



Gaussian Process

For any set \mathbf{S} , a Gaussian Process on \mathbf{S} is a set of random variables ($f(x), x \in \mathbf{S}$), so for any $n \in \mathbb{N}$ and $x_1, \dots, x_n \in \mathbf{S}$, $\{f(x_1), \dots, f(x_n)\}$ is (multivariate) Gaussian

- the number of elements in set \mathbf{S} can be infinite many.
- the number of random variables $f(x)$ can be any
 - $\{f(x_9), f(x_5), f(x_2), f(x_7)\}$ can build a mean and a covariance function, Gaussian!
 - $\{f(x_7), f(x_8)\}$ can build another mean and another covariance function, Gaussian!
- Gaussian Processes (GPs) are parameterised by a mean function $\mu(x)$, and a covariance function, or **kernel**, $K(x, x')$

$$\text{cov}(f_n, f_m) = \langle f_n - \mu(f_n), f_m - \mu(f_m) \rangle = \langle f_n, f_m \rangle - \mu(f_n)\mu(f_m)$$

Gaussian Process

For any set \mathbf{S} , a Gaussian Process on \mathbf{S} is a set of random variables $(f(x), x \in \mathbf{S})$, so for any $n \in \mathbb{N}$ and $x_1, \dots, x_n \in \mathbf{S}$, $\{f(x_1), \dots, f(x_n)\}$ is (multivariate) Gaussian

- the number of elements in set \mathbf{S} can be infinite many.
- the number of random variables $f(x)$ can be any
 - $\{f(x_9), f(x_5), f(x_2), f(x_7)\}$ can build a mean and a covariance function, Gaussian!
 - $\{f(x_7), f(x_8)\}$ can build another mean and another covariance function, Gaussian!
- Gaussian Processes (GPs) are parameterised by a mean function $\mu(x)$, and a covariance function, or **kernel**, $K(x, x')$

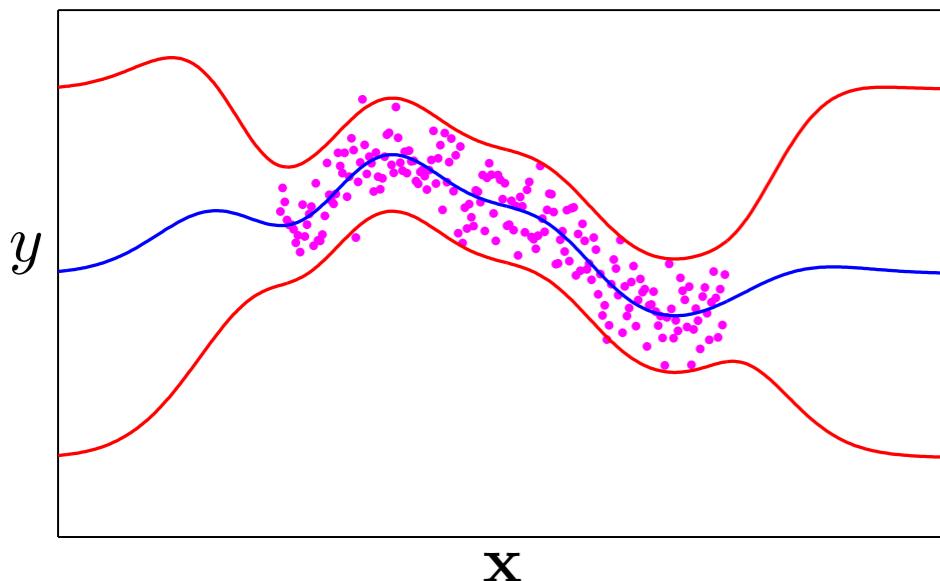


[any possible collection of random variables (from training data) creates a function]

☞ **Gaussian process is a distribution over functions**

Gaussian Process Regression

- Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^N\} = (\mathbf{X}, \mathbf{y})$
 - regression in not professional statement:
given a bunch of training data, and predict new y for a test input x
- We want to learn a function f with error bars from data \mathcal{D}



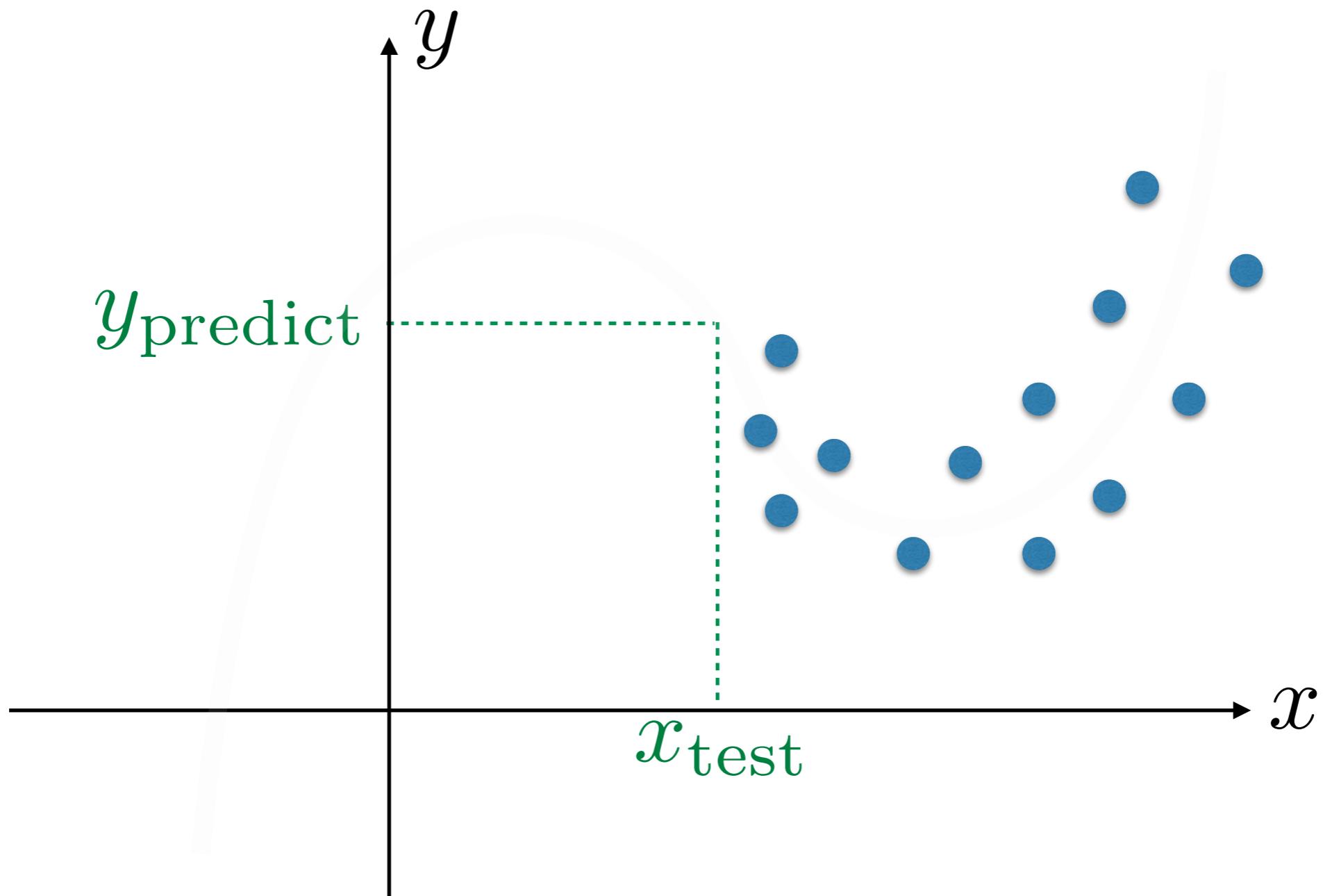
Gaussian process

$$y_n = f(\mathbf{x}_n) + \epsilon_n$$
$$\epsilon_n \sim \mathcal{N}(\cdot | 0, \beta^{-1})$$

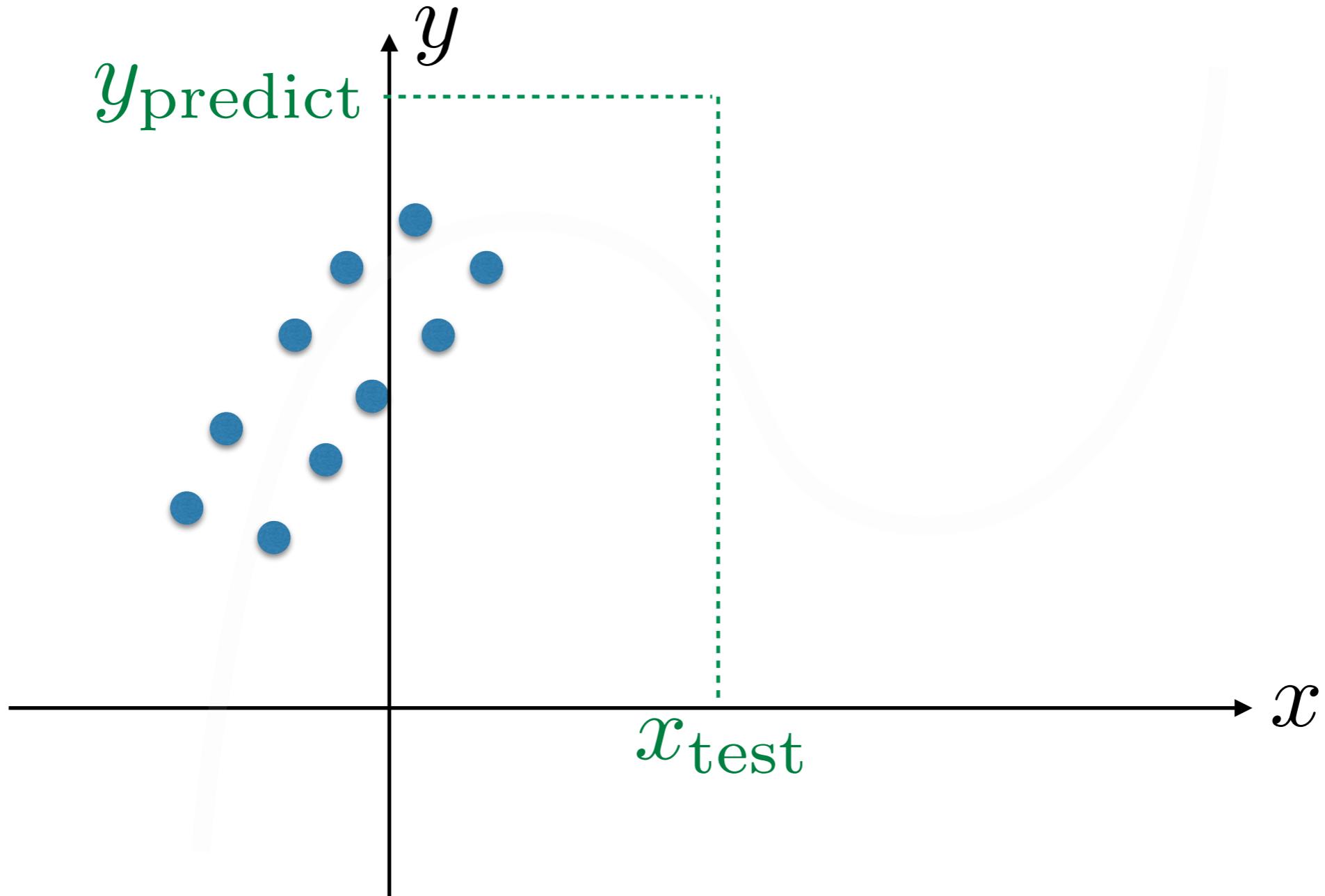
Gaussian Process Regression

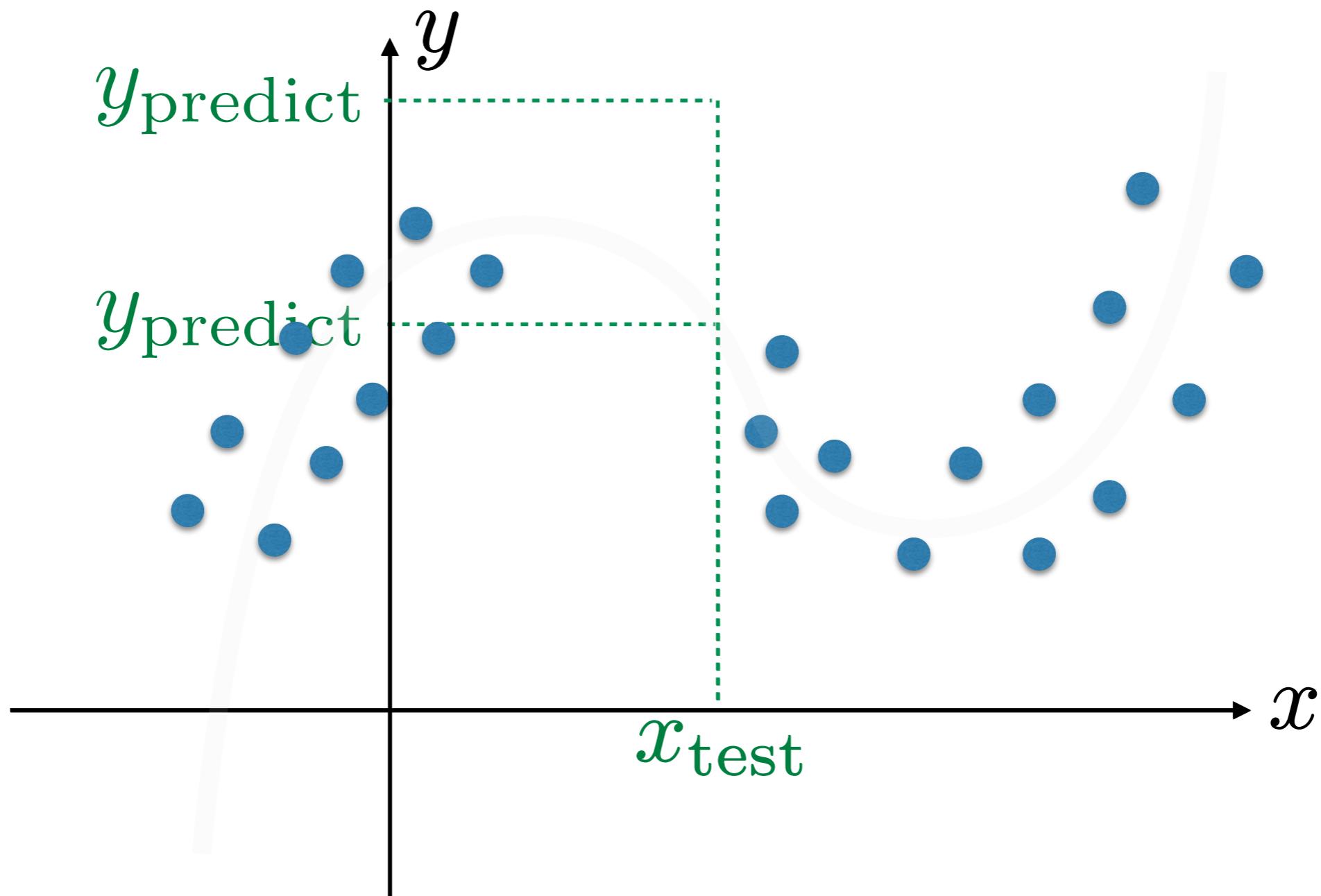
- Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^n\} = (X, \mathbf{y})$
 - regression is not professional statement:
given a bunch of training data, and predict new y for a test input \mathbf{x}
- Very basic ideas help you to understand GP regression:
 - It is a memory-based method!!!
 - Prediction is based on the relation between test input and training data!
 - Intuition: if \mathbf{x} and \mathbf{x}' are close to each other (in feature space),
then their y will be also close

Gaussian Process Regression



Gaussian Process Regression





Gaussian Process Regression

- Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^n\} = (X, \mathbf{y})$
 - regression is not professional statement:
given a bunch of training data, and predict new y for a test input \mathbf{x}
- Very basic ideas help you to understand GP regression:
 - It is a memory-based method!!!
 - Prediction is based on the relation between test input and training data!
 - Intuition: if \mathbf{x} and \mathbf{x}' are close to each other (in feature space),
then their y will be also close
 - metric to evaluate the similarity/distance between \mathbf{x} ↗ kernel!
 - GP is parameterised by a mean function $\mu(\mathbf{x})$, and a covariance
function, or kernel, $K(\mathbf{x}, \mathbf{x}')$.
 - Since random variable $y_n = f(\mathbf{x}_n)$, we can compute the similarity
between random variables by covariance

Gaussian Process Regression

- Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^n\} = (\mathbf{X}, \mathbf{y})$

regression is not professional statement:

The kernel function that determines K should be chosen to express the property that: for points \mathbf{x}_n and \mathbf{x}_m that are similar, the corresponding values $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ will be more strongly correlated than for dissimilar points.

- Intuition: if \mathbf{x} and \mathbf{x}' are close to each other (**in feature space**), then their y will be also close

- metric to evaluate the similarity/distance between \mathbf{x} ↩ kernel!

- GP is parameterised by a mean function $\mu(\mathbf{x})$, and a covariance function, or kernel, $K(\mathbf{x}, \mathbf{x}')$.

- Since random variable $y_n = f(\mathbf{x}_n)$, we can compute the similarity between random variables by covariance



ah-ha!

Gaussian Process Regression

- Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^n\} = (X, \mathbf{y})$
 - regression in not professional statement:
given a bunch of training data, and predict new y for a test input \mathbf{x}
- Very basic ideas help you to understand GP regression:

$$y_n = f(\mathbf{x}_n) + \epsilon_n$$
$$\epsilon_n \sim \mathcal{N}(\cdot | 0, \beta^{-1})$$

We don't have any parametric model like $y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$
Gaussian process regression is nonparametric!

Gaussian Process Regression

- Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^n\} = (X, \mathbf{y})$
 - regression in not professional statement:
given a bunch of training data, and predict new y for a test input \mathbf{x}
- Very basic ideas help you to understand GP regression:

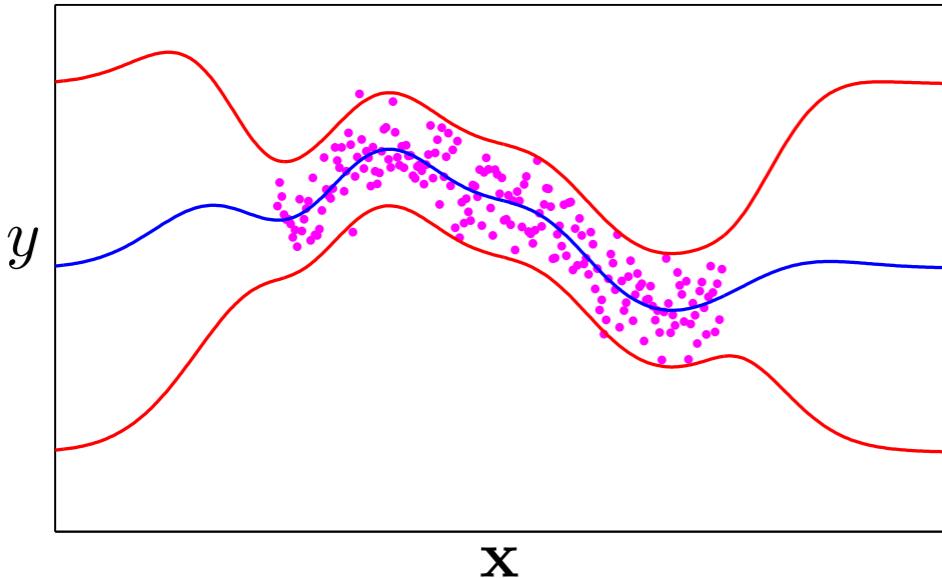
$$y_n = f(\mathbf{x}_n) + \epsilon_n$$
$$\epsilon_n \sim \mathcal{N}(\cdot | 0, \beta^{-1})$$

We don't have any parametric model like $y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$
Gaussian process regression is nonparametric!



“nonparametric” doesn’t mean there is no parameters, instead it means: the number of parameters will grows with number of data!
ah-ha!

Gaussian Process Regression



$$y_n = f(\mathbf{x}_n) + \epsilon_n$$
$$\epsilon_n \sim \mathcal{N}(\cdot | 0, \beta^{-1})$$

for $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top$ and $\mathbf{y} = [y_1, \dots, y_N]^\top$

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \beta^{-1}\mathbf{I}_N)$$

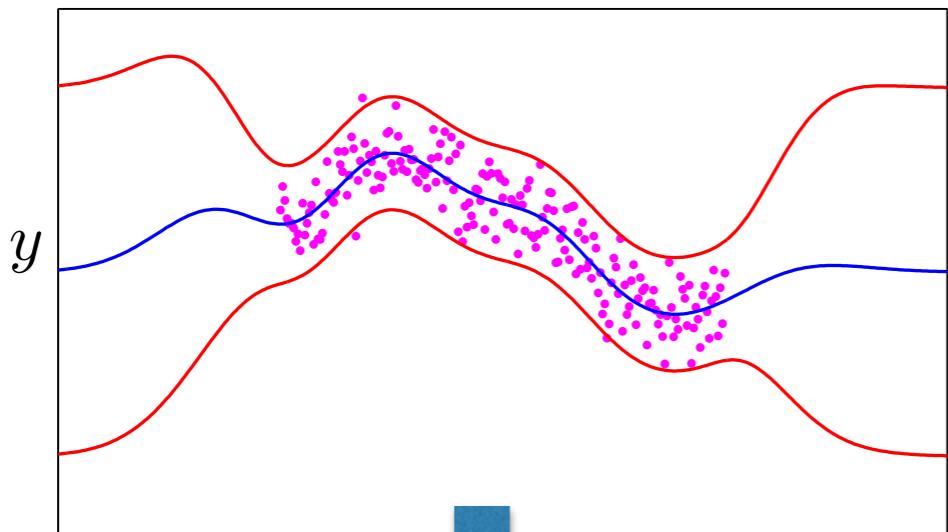
$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

marginal likelihood $p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C})$

where the covariance matrix \mathbf{C} has elements

$$\mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}$$

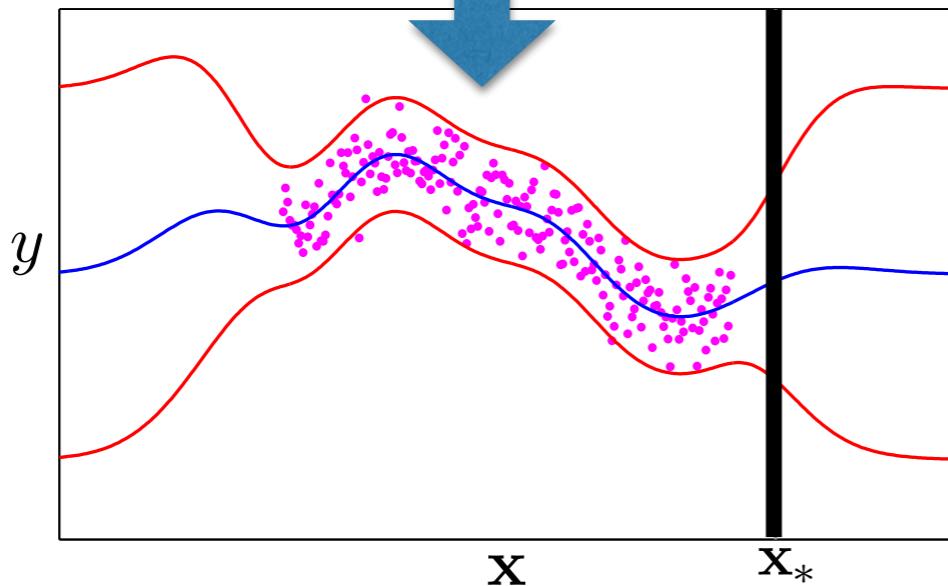
Gaussian Process Regression



marginal likelihood

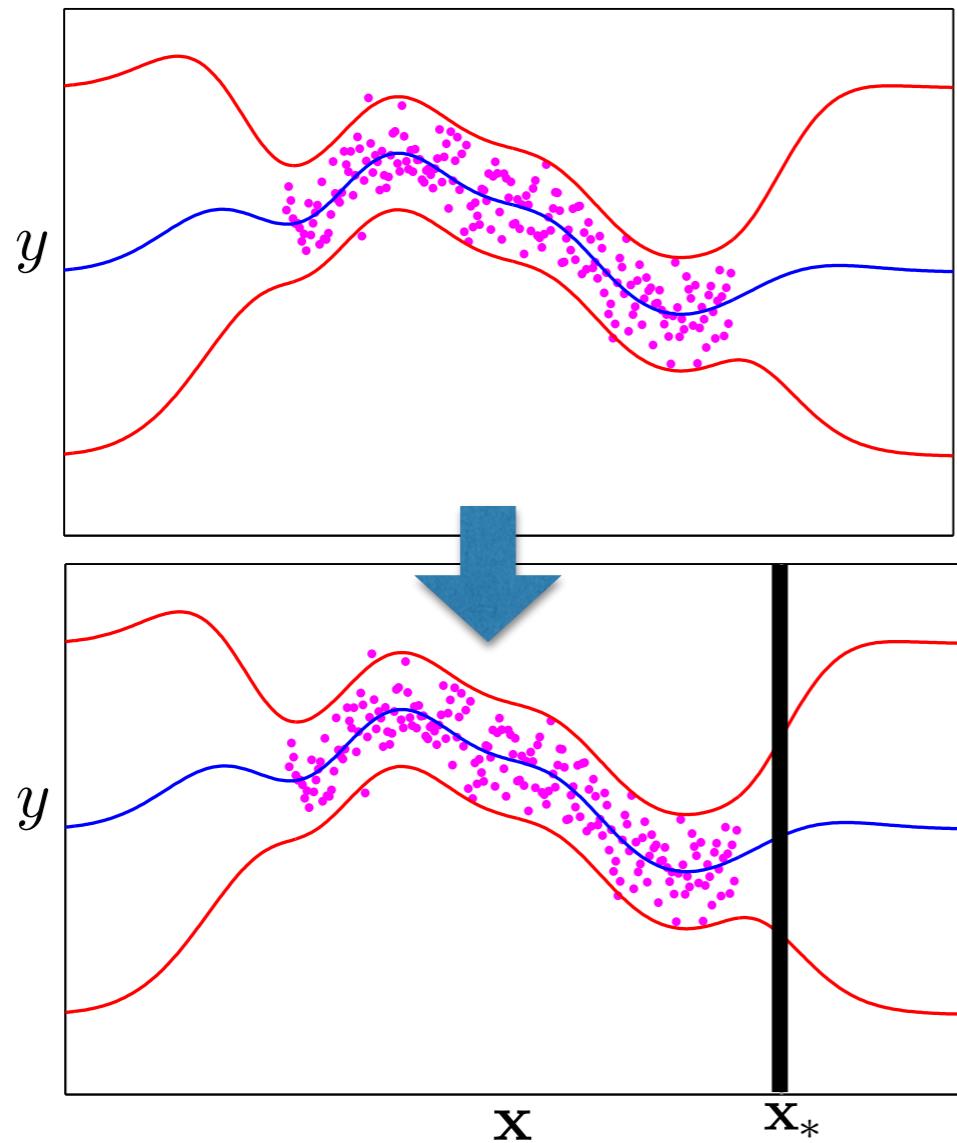
$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C})$$

$$\mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}$$



prediction?

Gaussian Process Regression



marginal likelihood

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C})$$

$$\mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}$$

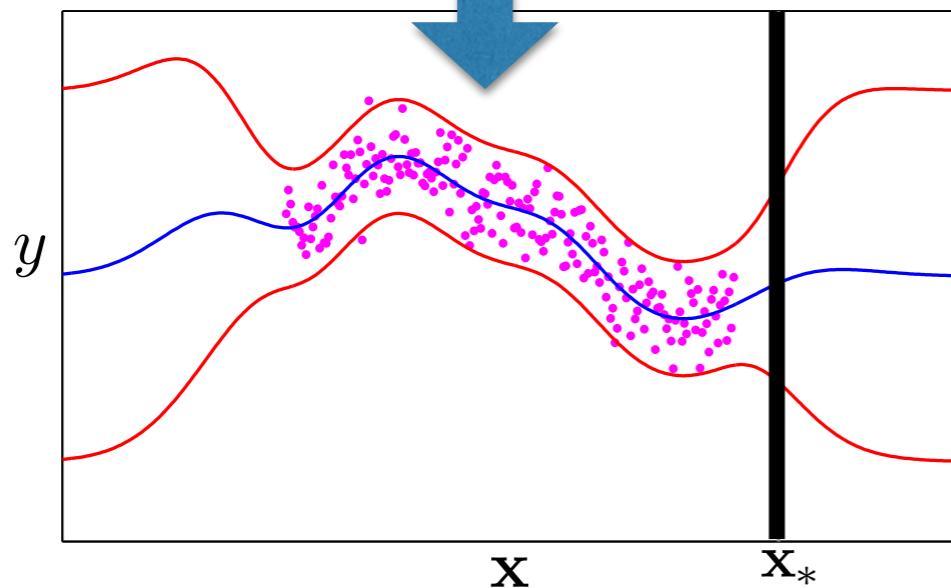
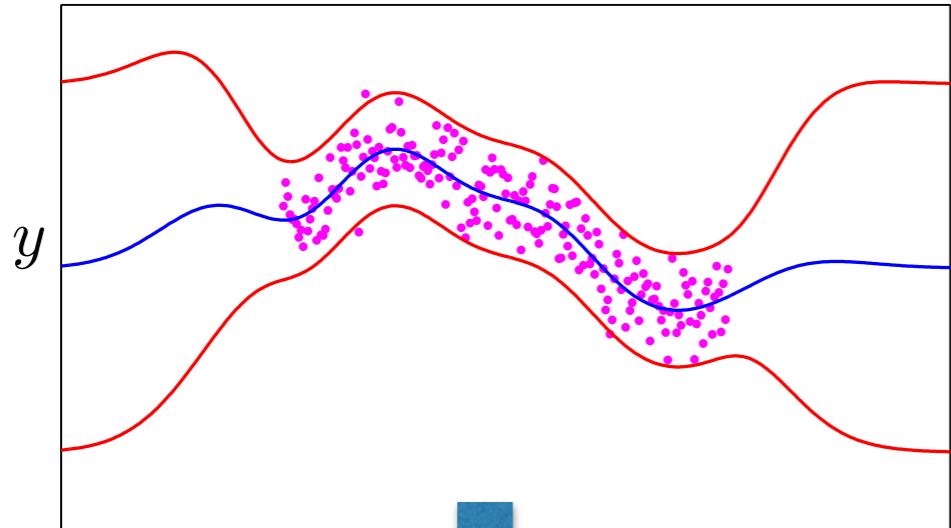
prediction?



ah-ha!

Because we know that all \mathbf{f} (training+testing) together become a multivariate Gaussian distribution \mathbf{G} , therefore if we want to predict the distribution of new \mathbf{f}^* , we just need to compute the covariance matrix of \mathbf{G} , then cut \mathbf{G} on \mathbf{f}^* to see the conditional distribution thus achieve prediction :D

Gaussian Process Regression



marginal likelihood

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C})$$

$$\mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}$$

prediction

denote $\mathbf{y}_{N+1} = [\mathbf{y}, y^*]^\top$ and $y^* = f(\mathbf{x}^*)$

$$p(\mathbf{y}_{N+1}) = \mathcal{N}(\mathbf{y}_{N+1}, |\mathbf{0}, \mathbf{C}_{N+1})$$

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C} & k(\mathbf{x}, \mathbf{x}^*) \\ k(\mathbf{x}, \mathbf{x}^*)^\top & k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1} \end{bmatrix}$$

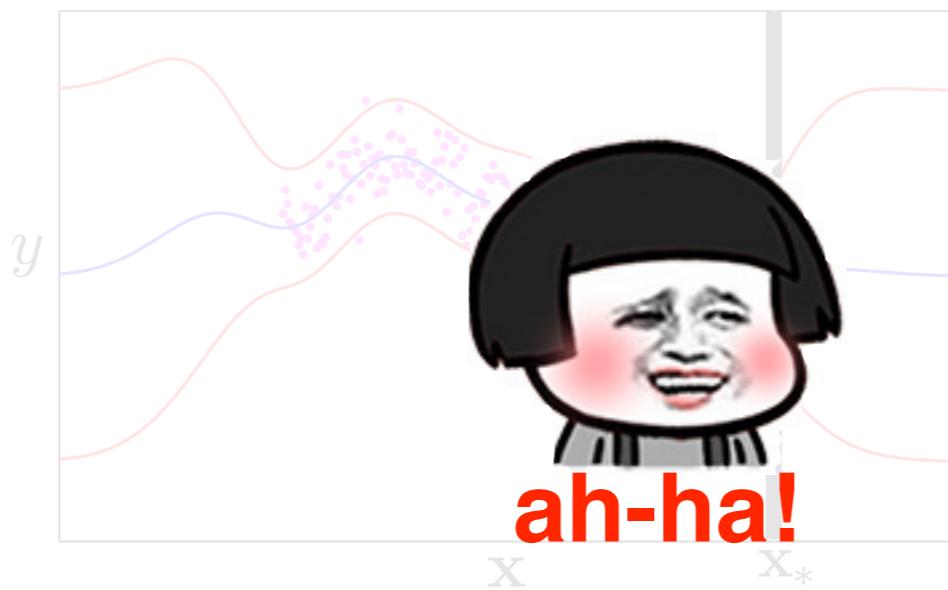
conditional distribution $p(y^*|\mathbf{y})$ is a Gaussian distribution with:

$$\mu(\mathbf{x}^*) = k(\mathbf{x}, \mathbf{x}^*)^\top \mathbf{C}^{-1} \mathbf{y}$$

$$\sigma^2(\mathbf{x}^*) = k^* - k(\mathbf{x}, \mathbf{x}^*)^\top \mathbf{C}^{-1} k(\mathbf{x}, \mathbf{x}^*)$$

$$k^* = k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1}$$

Gaussian Process Regression



prediction

denote $\mathbf{y}_{N+1} = [\mathbf{y}, y^*]^\top$ and $y^* = f(\mathbf{x}^*)$
 $p(\mathbf{y}_{N+1}) = \mathcal{N}(\mathbf{y}_{N+1}, \mathbf{0}, \mathbf{C}_{N+1})$

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C} & k(\mathbf{x}, \mathbf{x}^*) \\ k(\mathbf{x}, \mathbf{x}^*)^\top & k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1} \end{bmatrix}$$

1° kernel

conditional distribution $p(y^* | \mathbf{y})$ is a Gaussian distribution with:

2° conditional

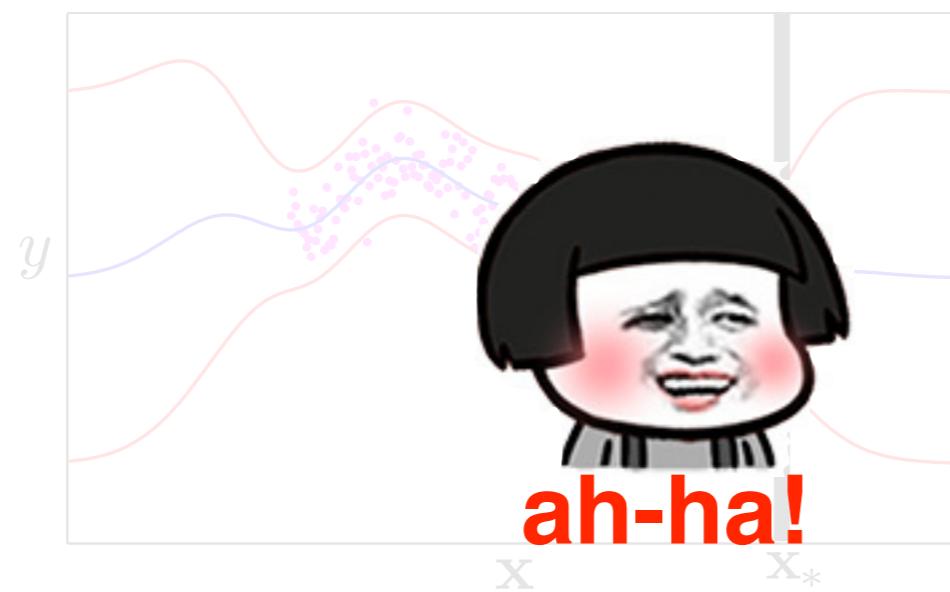
$$\mu(\mathbf{x}^*) = k(\mathbf{x}, \mathbf{x}^*)^\top \mathbf{C}^{-1} \mathbf{y}$$

3° done!

$$\sigma^2(\mathbf{x}^*) = k^* - k(\mathbf{x}, \mathbf{x}^*)^\top \mathbf{C}^{-1} k^*$$

$$k^* = k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1}$$

Gaussian Process Regression



similarity from testing \mathbf{x}^*
to all the training data!
 $\mathbf{C}^{-1}\mathbf{y}$ is just a vector of scalars

$$p(\mathbf{y}_{N+1}) = \mathcal{N}(\mathbf{y}_{N+1}, \mathbf{0}, \mathbf{C}_{N+1})$$

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C} & k(\mathbf{x}, \mathbf{x}^*) \\ k(\mathbf{x}, \mathbf{x}^*)^\top & k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1} \end{bmatrix}$$

conditional distribution $p(\mathbf{y}^* | \mathbf{y})$ is a Gaussian distribution with:

2° conditional

$$\mu(\mathbf{x}^*) = k(\mathbf{x}, \mathbf{x}^*)^\top \mathbf{C}^{-1} \mathbf{y}$$

3° done!

$$\sigma^2(\mathbf{x}^*) = k^* - k(\mathbf{x}, \mathbf{x}^*)^\top \mathbf{C}^{-1} k^*$$

$$k^* = k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1}$$

Gaussian Process: learning the kernel

Still remember this?

The kernel function that determines K should be chosen to express the property that:
for points \mathbf{x}_n and \mathbf{x}_m that are similar,
the corresponding values $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ will be more strongly correlated than for dissimilar points.

- Consider covariance function \mathbf{C} with hyper-parameters θ

$$k_\theta(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left\{-\theta_1 \frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{2}\right\} + \theta_2 + \theta_3 \mathbf{x}_n^\top \mathbf{x}_m$$

Gaussian Process: learning the kernel

Still remember this?

The kernel function that determines K should be chosen to express the property that:
for points \mathbf{x}_n and \mathbf{x}_m that are similar,
the corresponding values $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ will be more strongly correlated than for dissimilar points.

- Consider covariance function \mathbf{C} with hyper-parameters $\boldsymbol{\theta}$

$$k_{\theta}(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left\{-\theta_1 \frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{2}\right\} + \theta_2 + \theta_3 \mathbf{x}_n^\top \mathbf{x}_m$$

- Given $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^N\} = (\mathbf{X}, \mathbf{y})$, the marginal likelihood is function of $\boldsymbol{\theta}$

$$p(\mathbf{y}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_{\boldsymbol{\theta}})$$

$$\ln p(\mathbf{y}|\boldsymbol{\theta}) = -\frac{1}{2} \ln |\mathbf{C}_{\boldsymbol{\theta}}| - \frac{1}{2} \mathbf{y}^\top \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{y} - \frac{N}{2} \ln (2\pi) \quad \text{👉} \quad \frac{\partial \ln p(\mathbf{y}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

Gaussian Process: learning the kernel

Still remember this?

The kernel function that determines K should be chosen to express the property that:
for points x_n and x_m that are similar,
the corresponding values $y(x_n)$ and $y(x_m)$ will be more strongly correlated than for dissimilar points.

- Consider covariance function C with hyper-parameters θ

The keys in Gaussian Process Regression are:

1. Choose kernel function
2. Estimate the proper hyper-parameters!

$$p(\mathbf{y}|\theta) = \mathcal{N}(\mathbf{y}|0, \mathbf{C}_\theta)$$

$$\ln p(\mathbf{y}|\theta) = -\frac{1}{2} \ln |\mathbf{C}_\theta| - \frac{1}{2} \mathbf{y}^\top \mathbf{C}_\theta^{-1} \mathbf{y} - \frac{N}{2} \ln (2\pi)$$

$$\frac{\partial \ln p(\mathbf{y}|\theta)}{\partial \theta}$$

Feature selection can be achieved by ARD in GPs

Problem: Often there are *many* possible inputs that might be relevant to predicting a particular output. We need algorithms that automatically decide which inputs are relevant.

Automatic Relevance Determination:

Consider this covariance function:

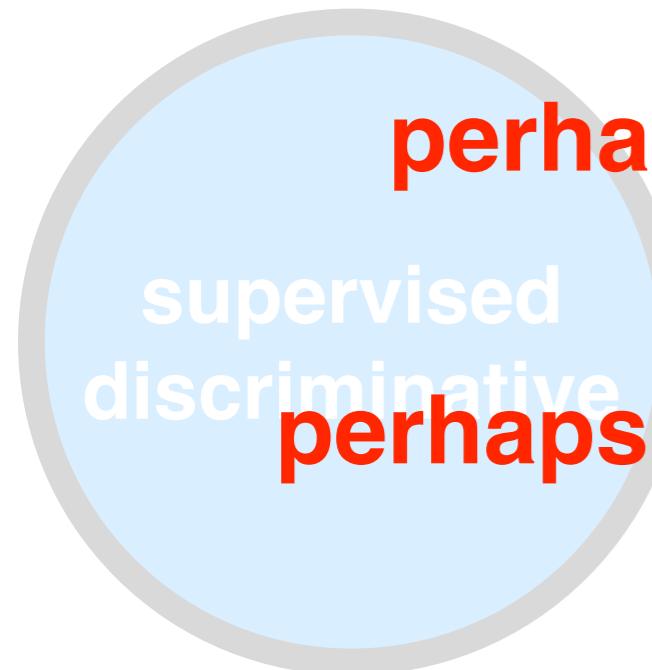
$$\mathbf{K}_{nn'} = v \exp \left[-\frac{1}{2} \sum_{d=1}^D \left(\frac{x_n^{(d)} - x_{n'}^{(d)}}{r_d} \right)^2 \right]$$

The parameter r_d is the length scale of the function along input dimension d .

As $r_d \rightarrow \infty$ the function f varies less and less as a function of $x^{(d)}$, that is, the d th dimension becomes *irrelevant*.

Given data, by learning the lengthscales (r_1, \dots, r_D) it is possible to do automatic feature selection.

Today



**perhaps most powerful regressor
before deep learning**

**perhaps most powerful classifier
before deep learning**



Dimension Reduction

Clustering

HMM

Graphical Models

Kernel
Methods

Gaussian
Process

Support Vector
Machine

Sparse kernel machines

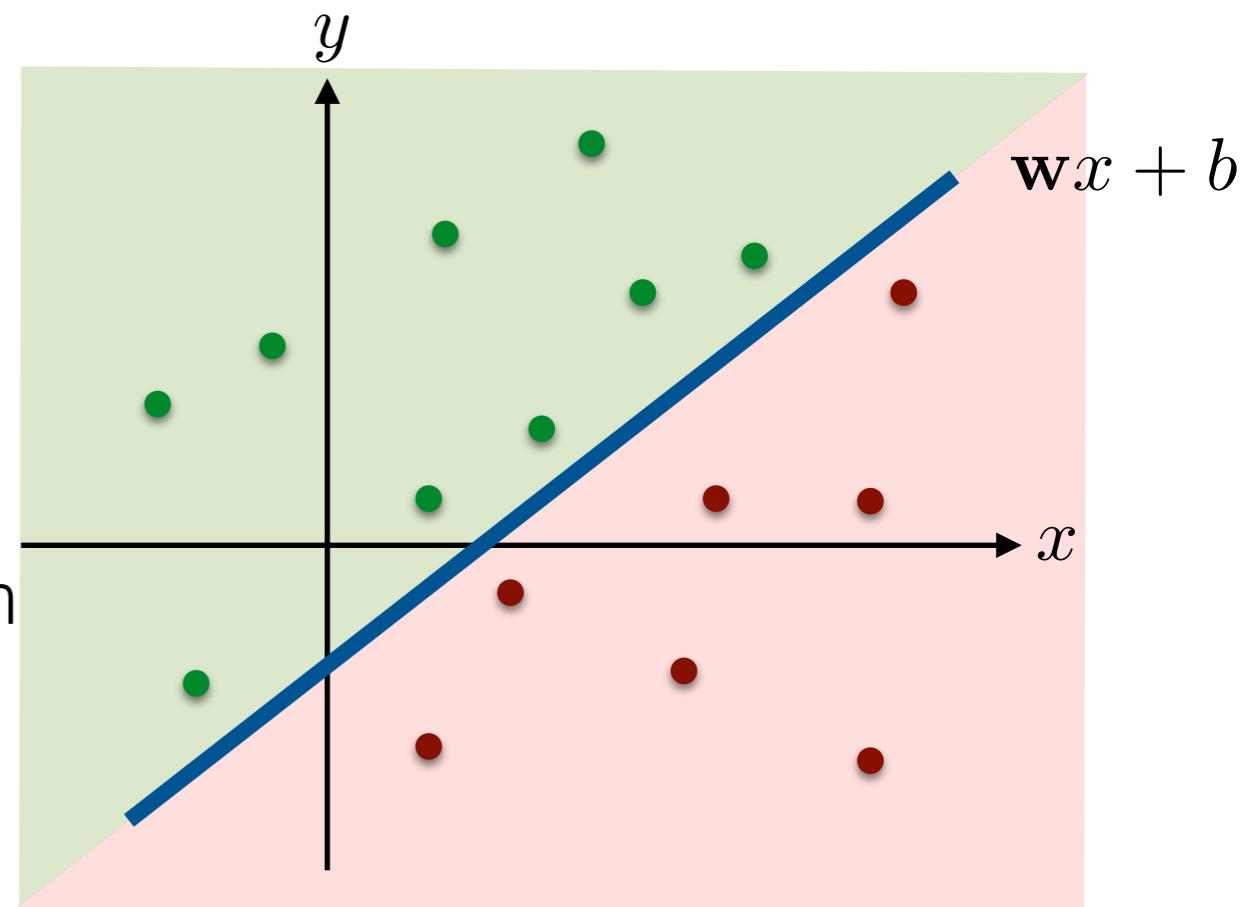
- Memory based methods: we would like to utilise the training data in the prediction!
 - What we have introduced (KNN, Parzen probability estimation, GPs) need to keep “ALL” training data ... heavy memory demands!
 - Support Vector Machine (SVM), what we are going to introduce, is a kernel-based algorithm that have **sparse** solutions!
 - Predictions by kernel functions evaluated on **subset** of training data!

Sparse kernel machines

- Let's start from simple **binary classification** problem:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

- classification based on sign of $y(\mathbf{x})$
- please note that explicit computing on $\phi(\mathbf{x})$ will be avoided by dual form
👉 kernel!

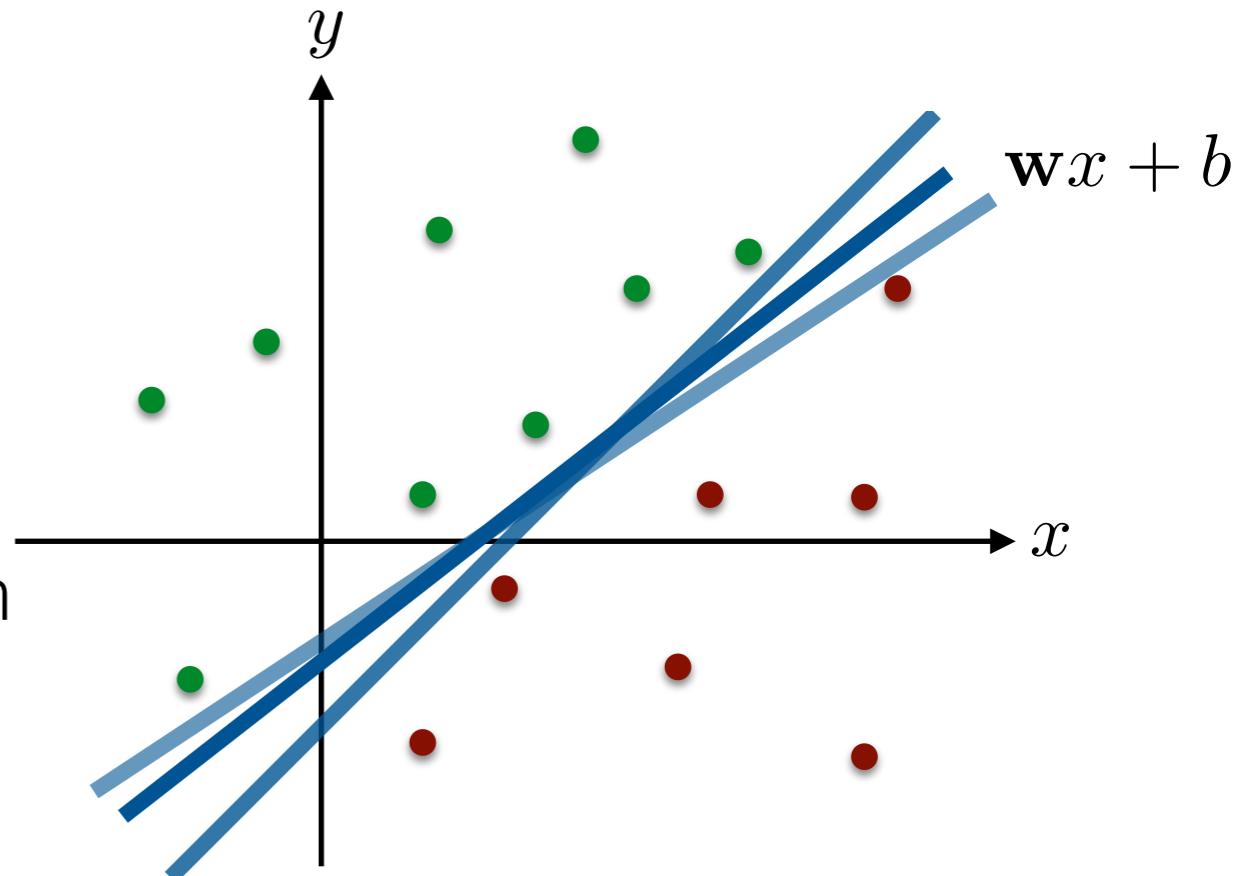


Sparse kernel machines

- Let's start from simple **binary classification** problem:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

- classification based on sign of $y(\mathbf{x})$
- please note that explicit computing on $\phi(\mathbf{x})$ will be avoided by dual form
👉 kernel!



- When assuming the data is linearly separable, it means there will be infinite solutions for \mathbf{w} and b , we should try to find the one with minimum generalisation error.



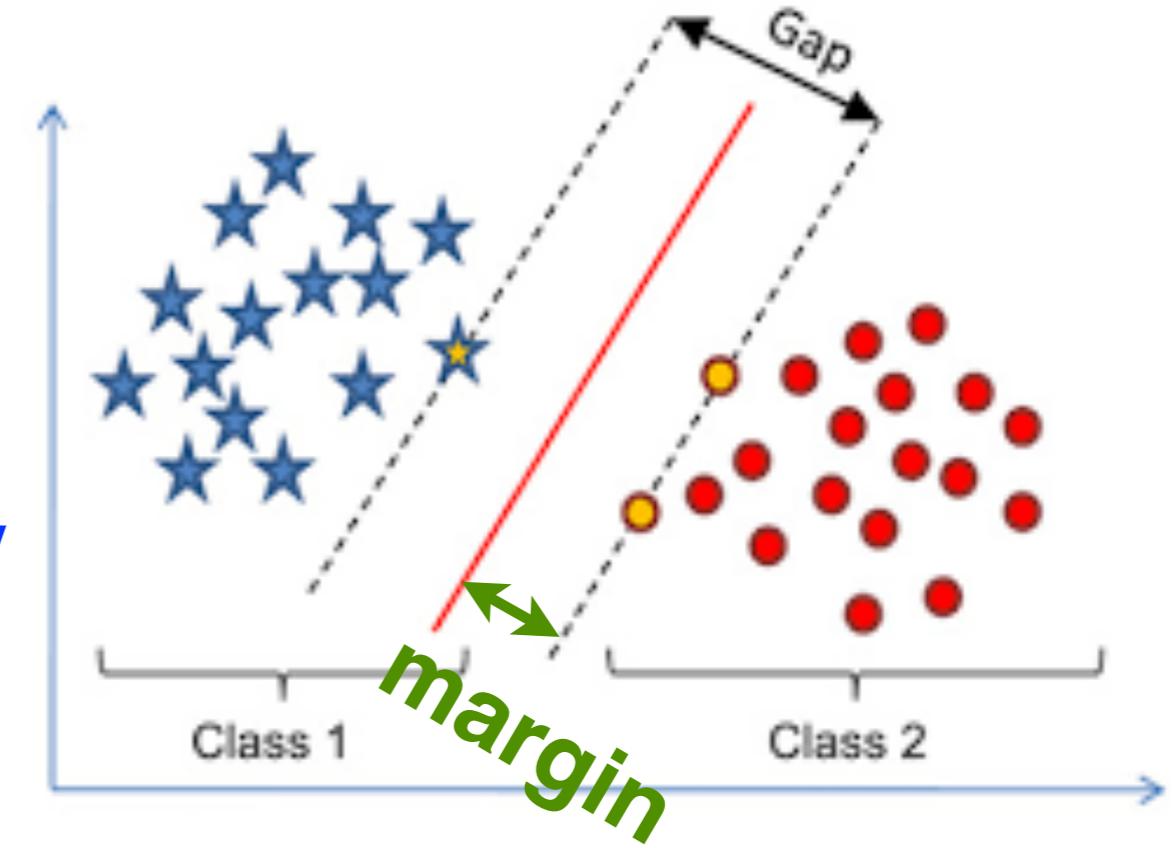
what does it mean?

Sparse kernel machines

- Let's start from simple **binary classification** problem:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

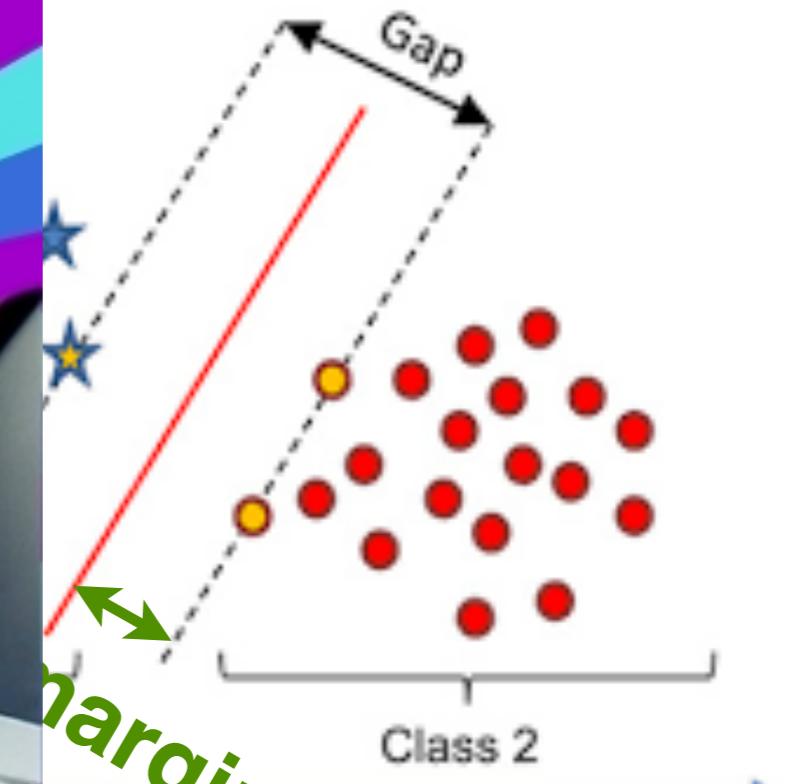
Maximise margin!
margin is the smallest distance b/w decision boundary and any of the training samples



- When assuming the data is linearly separable, it means there will be infinite solutions for \mathbf{w} and b , we should try to find the one with minimum generalisation error.

Sparse kernel machines

Maximise margin!
margin is the small
distance b/w decision
and any of the training



$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

perpendicular distance from a data point x to a hyperplane $y(\mathbf{x}) = 0$ is given by

$$\frac{|y(\mathbf{x})|}{\|\mathbf{w}\|}$$

Sparse kernel machines

we are only interested in solutions for which all data points are correctly classified

so that $t_n y(\mathbf{x}_n) > 0 \quad \forall n$

**t_n and $y(\mathbf{x}_n)$ with
the same sign**

$$\frac{t_n y(\mathbf{x})}{\|\mathbf{x}\|} = \frac{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{x}\|}$$

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

perpendicular distance from a data point x to a hyperplane $y(\mathbf{x}) = 0$ is given by

$$\frac{|y(\mathbf{x})|}{\|\mathbf{w}\|} \text{ absolute value}$$

Sparse kernel machines

we are only interested in solutions for which all data points are correctly classified so that $t_n y(\mathbf{x}_n) > 0 \quad \forall n$

$$\frac{t_n y(\mathbf{x})}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

we want to maximise the margin!

$$\underset{\mathbf{w}, b}{\operatorname{argmax}} \left\{ \min_n \frac{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right\}$$

margin is the smallest distance b/w decision boundary and any of the training samples

Sparse kernel machines

we are only interested in solutions for which all data points are correctly classified so that $t_n y(\mathbf{x}_n) > 0 \quad \forall n$

$$\frac{t_n y(\mathbf{x})}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

we want to maximise the margin!

$$\operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_n \frac{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right\}$$

note here if we scale \mathbf{w}, b together $\frac{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$ won't change

so we set a condition to get the unique solution:

for the points closest to the hyperplane, $t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b) = 1$
which means, for all data points, $t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1$

Sparse kernel machines

- Add these two criteria together!

$$\operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_n \frac{t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right\}$$
$$t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) = 1$$

- The overall optimisation problem becomes:

$$\text{maximise } \frac{1}{\|\mathbf{w}\|} \Rightarrow \text{minimise } \|\mathbf{w}\|$$

$$\text{subject to } t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1$$

- Here comes Lagrange multipliers $a_n \geq 0$ $\mathbf{a} = (a_1, \dots, a_N)^\top$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N a_n \{t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) - 1\}$$

remember regularisation story?

- minimise w.r.t. \mathbf{w} and b while maximise w.r.t. \mathbf{a}

Sparse kernel machines

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b) - 1\}$$

$$\begin{aligned}\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \mathbf{w}} &= 0 \quad \text{⇒} \quad \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \\ \frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial b} &= 0 \quad \text{⇒} \quad 0 = \sum_{n=1}^N a_n t_n\end{aligned}$$



ah-ha!

w is the linear combination based on
the feature representations of training data.
If you put this w back to L(w, b ,a),
it is easier to see that L has all x appearing as
scalar product! kernel trick!

Sparse kernel machines

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N a_n \{ t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b) - 1 \}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \mathbf{w}} = 0 \quad \text{shake hand} \quad \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad \text{cancel out}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial b} = 0 \quad \text{shake hand} \quad 0 = \sum_{n=1}^N a_n t_n$$

- put \mathbf{w} back to $L(\mathbf{w}, b, \mathbf{a})$ to get the **dual representation of the max-margin problem**, in which we maximise:

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0, n = 1, \dots, N$

$$\sum_{n=1}^N a_n t_n = 0$$

Sparse kernel machines

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0, n = 1, \dots, N$

$$\sum_{n=1}^N a_n t_n = 0$$

use SMO (Sequential Minimal optimisation) for QP (Quadratic Programming) to solution of a objective is quadratic and so any local optimum will also be a global optimum provided that the constraints define a convex region (since being linear)

Sparse kernel machines

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0, n = 1, \dots, N$

$$\sum_{n=1}^N a_n t_n = 0$$

after getting the solution for \mathbf{a} , the prediction can be done by:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

☞ $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n \underline{k(\mathbf{x}, \mathbf{x}_n)} + b$

kernel!

memory-based method!

Support Vector Machine

- See from **KKT condition**

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^\top \phi(\mathbf{x}_n) + b) - 1\}$$

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0$$

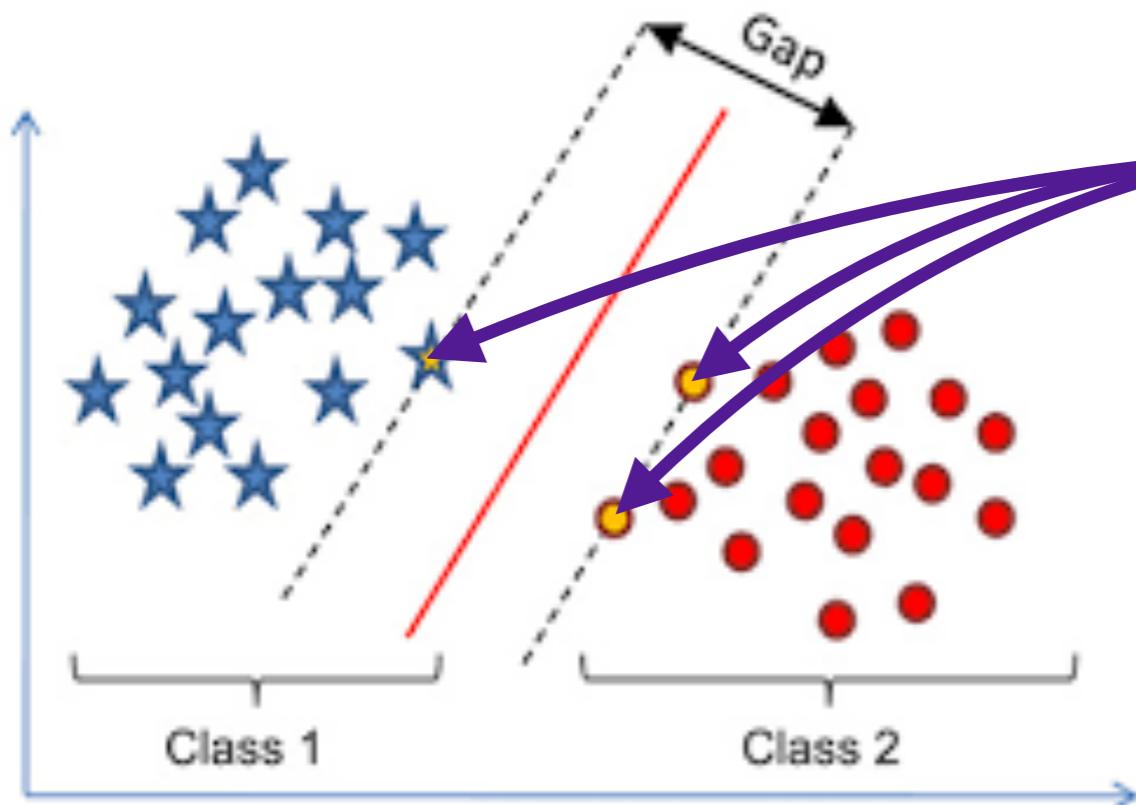
for every data point, either $a_n = 0$ or $t_n y(\mathbf{x}_n) = 1$

**data points having this condition
won't be used for prediction**

**data points having this condition
are called: support vector**

prediction: $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$

Support Vector Machine



These data points (vectors) support the hyperplane! we can call them support vectors!

$$t_n y(\mathbf{x}_n) = 1$$

data points having this condition are called: **support vector**

remember few slides ago?

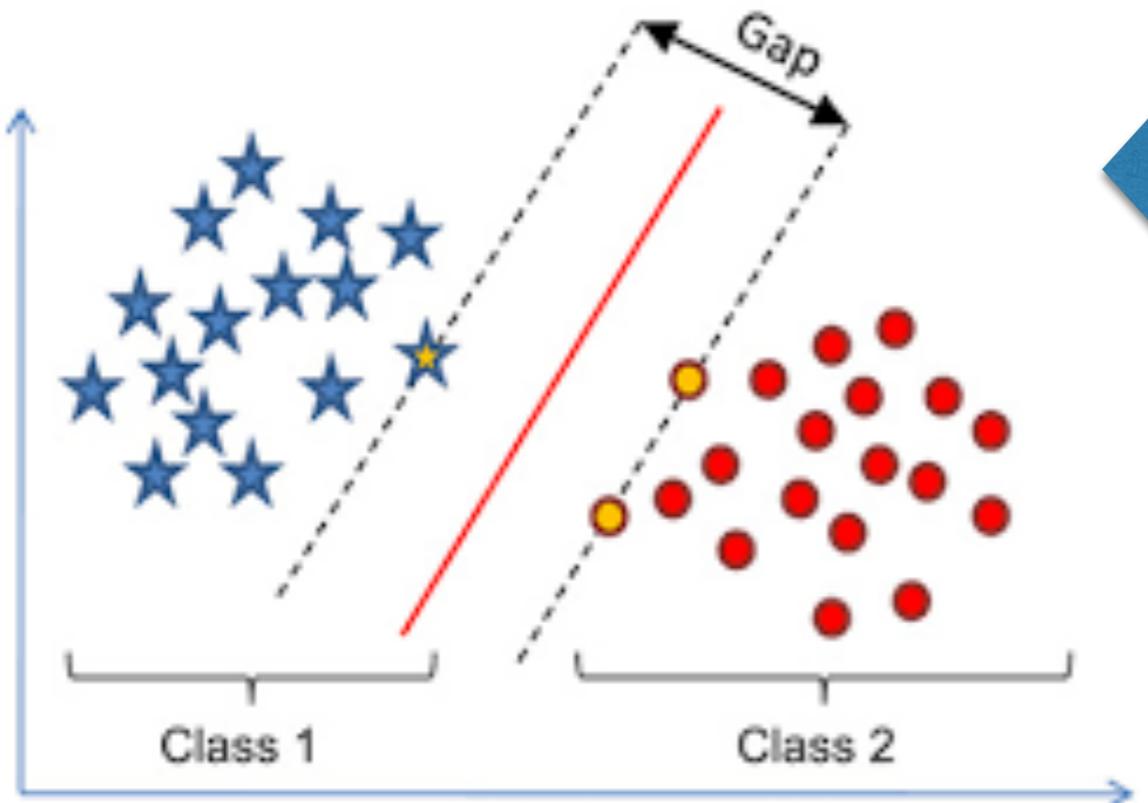
Note here if we scale w, σ together $\frac{t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$ won't change

so we set a condition to get the unique solution:

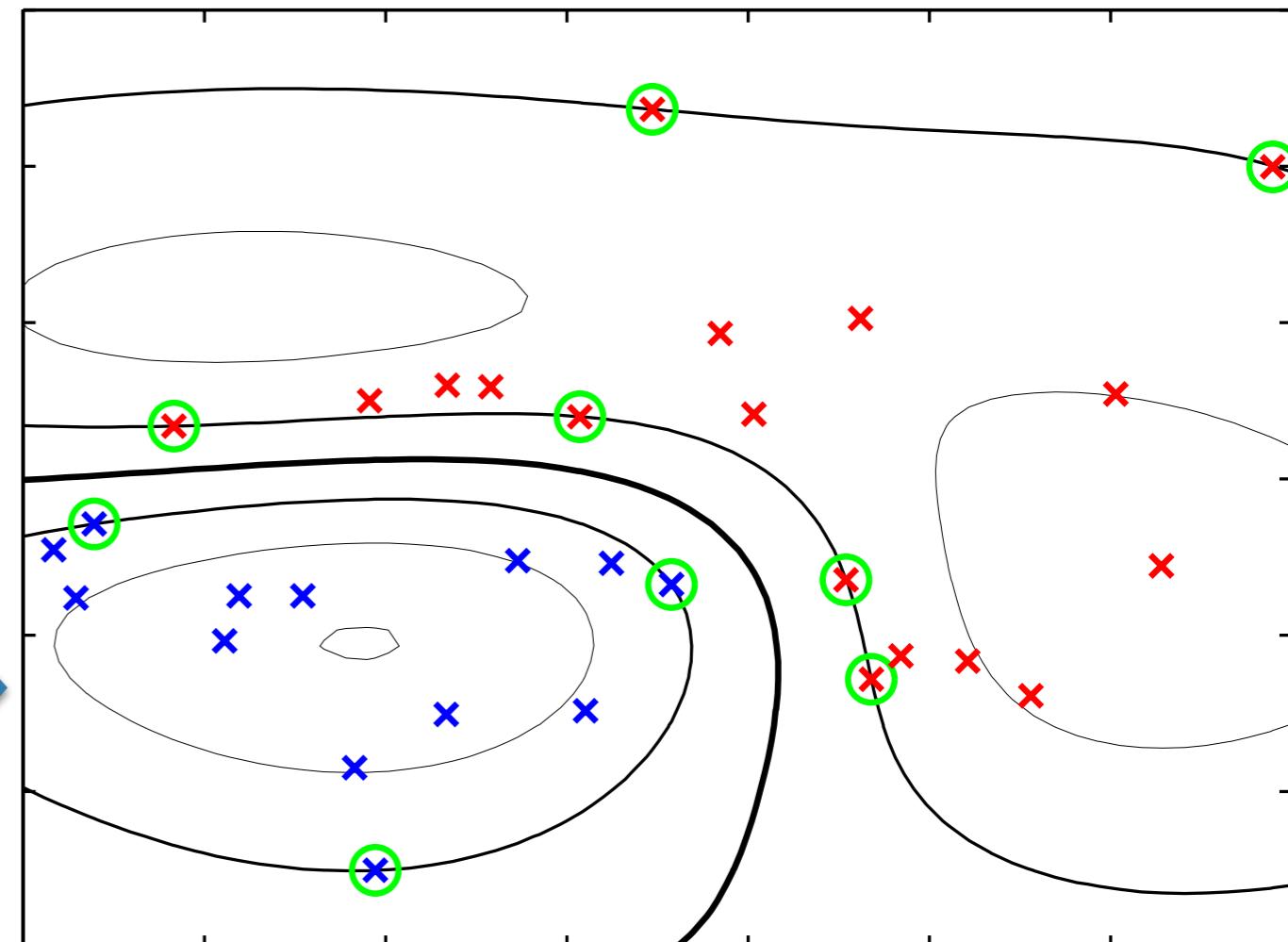
for the points closest to the hyperplane, $t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) = 1$

which means, for all data points, $t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1$

Support Vector Machine



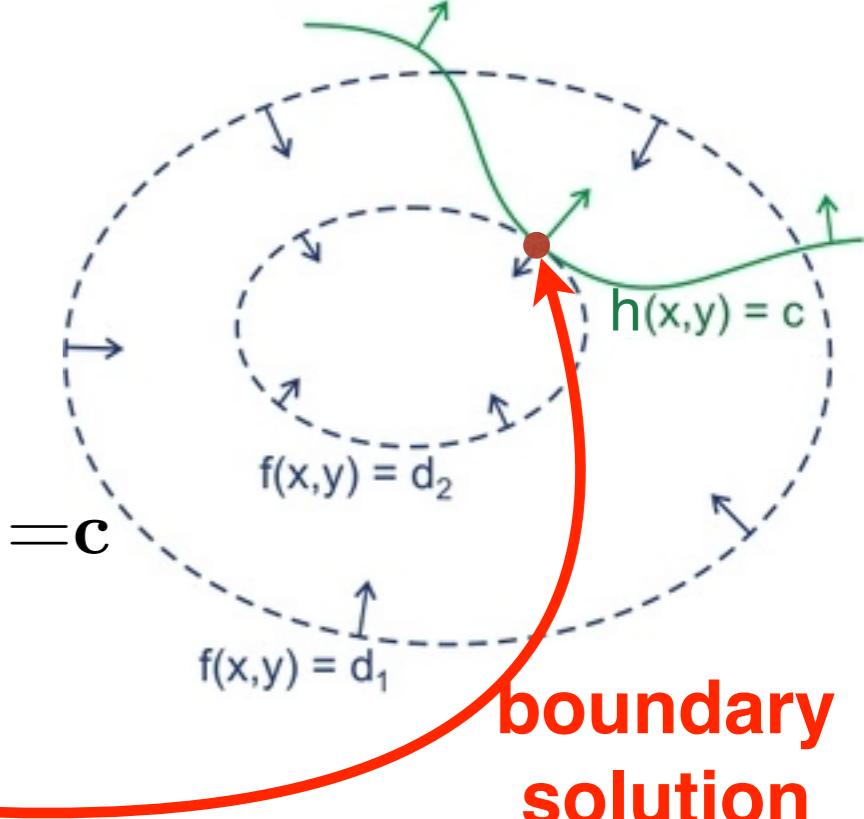
linearly separable in feature space
(whose dimensionality could be
much higher due to kernel function)



how decision boundary could
be in original data space

KKT condition

- simple optimisation problem
 - Green line represents equality constraint $h(x,y)=c$
 - Blue line shows contour lines of function f
 - Intuitively you can see in the optimum, the gradient of f is parallel to the gradient of h



**boundary
solution**

$$\text{let } L(X, \lambda, \mu) = f(X) + \sum_{j=1}^p \lambda_j h_j(X) + \sum_{k=1}^q \mu_k g_k(X)$$

$$\nabla f \in \text{span}\{\nabla h\}$$

$$\text{i.e. } \nabla f = -\eta \nabla h$$

equality inequality
constraints constraints

KKT condition

$$\left. \frac{\partial L}{\partial X} \right|_{X=X^*} = 0$$

equality constraints should be satisfied

$$\lambda_j \neq 0,$$

$$u_k \geq 0,$$

$$u_k g_k(X^*) = 0$$

$$h_j(X^*) = 0 \quad j=1,2,\dots,p$$

$$g_k(X^*) \leq 0 \quad k=1,2,\dots,q$$

This holds for both interior or boundary solution, called as complementary slackness.
(for interior solution, inequality constraints inactive, $u=0$)

inequality constraints to be active

Support Vector Machine - with Soft Margin!

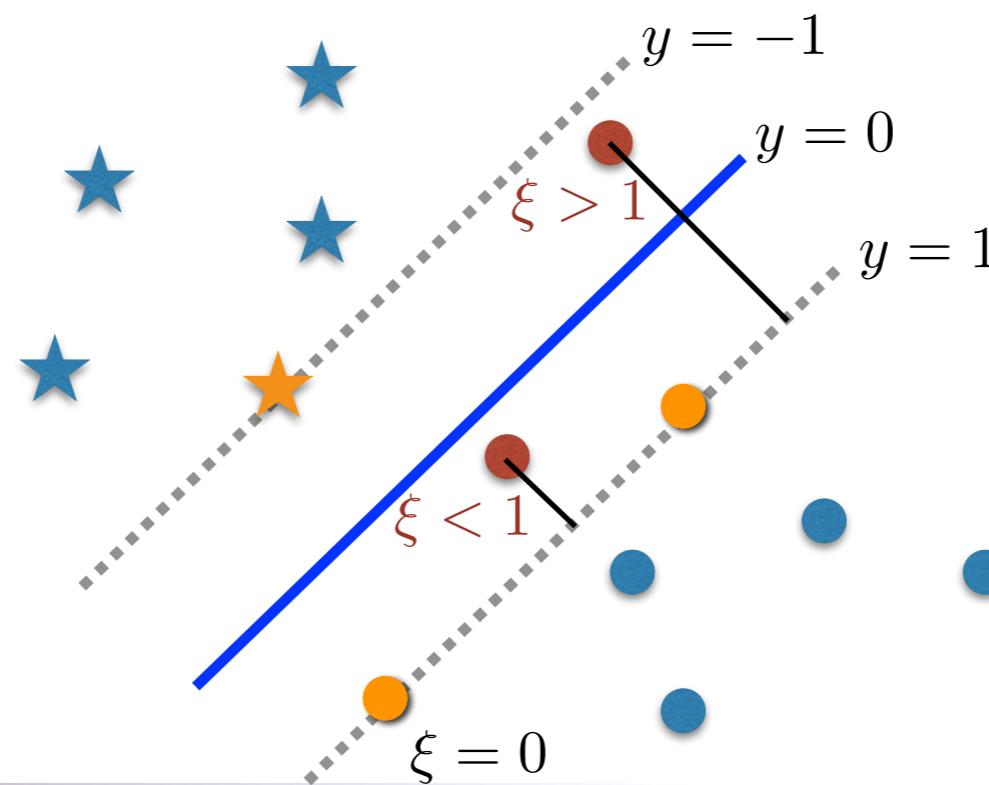
- Data points are allowed to be on the “wrong side” of boundary
 - But with a penalty ~~inversely~~ proportional to the distance from boundary

slack: $\xi_n = |t_n - y(\mathbf{x}_n)|$

- inequality $t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1$ replaced by $t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n$
also $\xi_n \geq 0$
- optimisation problem now as
subject to two constraints above

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

also minimise slack



Support Vector Machine - with Soft Margin!

- Data points are allowed to be on the “wrong side” of boundary
 - But with a penalty inversely proportional to the distance from boundary

slack: $\xi_n = |t_n - y(\mathbf{x}_n)|$

- inequality $t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1$ replaced by $t_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n$
also $\xi_n \geq 0$
- optimisation problem now as
subject to two constraints above

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

also minimise slack

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

- KKT condition

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 + \xi_n \geq 0$$

$$a_n (t_n y(\mathbf{x}_n) - 1 + \xi_n) = 0$$

$$\mu_n \geq 0$$

$$\xi_n \geq 0$$

$$\mu_n \xi_n = 0$$

Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

- KKT condition

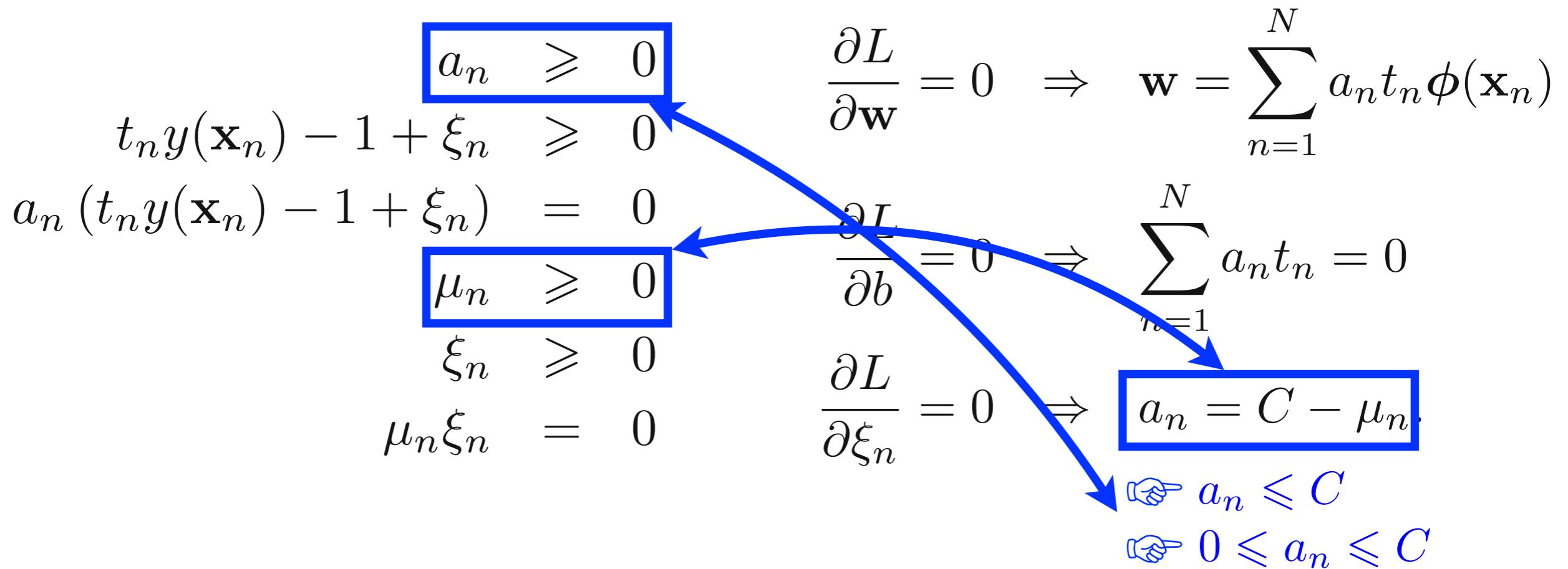
$$\begin{array}{lcl} a_n & \geqslant & 0 \\ t_n y(\mathbf{x}_n) - 1 + \xi_n & \geqslant & 0 \\ a_n (t_n y(\mathbf{x}_n) - 1 + \xi_n) & = & 0 \\ \mu_n & \geqslant & 0 \\ \xi_n & \geqslant & 0 \\ \mu_n \xi_n & = & 0 \end{array} \quad \begin{array}{l} \frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \\ \frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n t_n = 0 \\ \frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n. \end{array}$$

Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

- KKT condition



Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

put back to $L(\mathbf{w}, b, \mathbf{a})$

👉 $\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$

subject to $0 \leq a_n \leq C$

$$\sum_{n=1}^N a_n t_n = 0$$

use SMO for QP to get solution of a

Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

likewise, after getting the solution for \mathbf{a} , the prediction can be done by:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$
$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$



$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n \underline{k(\mathbf{x}, \mathbf{x}_n)} + b$$

kernel!

memory-based method!

Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

- KKT condition

$$\begin{array}{rcl} a_n & \geqslant & 0 \\ t_n y(\mathbf{x}_n) - 1 + \xi_n & \geqslant & 0 \\ a_n (t_n y(\mathbf{x}_n) - 1 + \xi_n) & = & 0 \end{array}$$



for every data point,
either $a_n = 0$ or $t_n y(\mathbf{x}_n) = 1 - \xi_n$
not used in prediction **data points related to support vector**

$$\mu_n \geqslant 0$$

$$\xi_n \geqslant 0$$

$$\mu_n \xi_n = 0$$

prediction: $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$

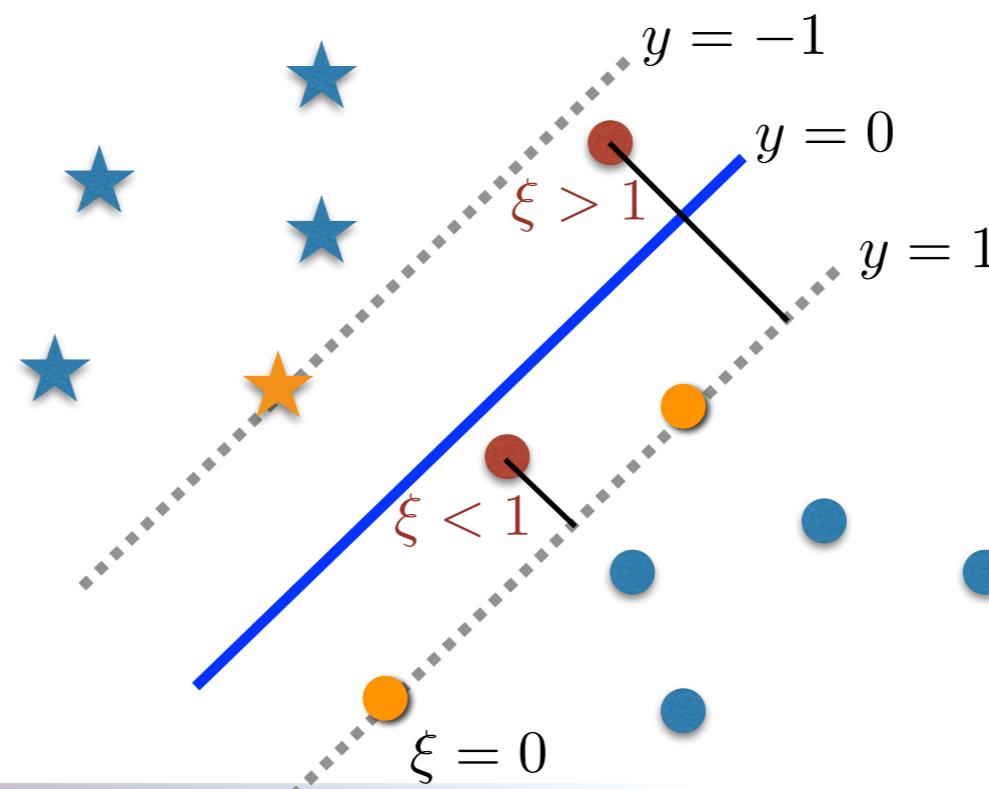
Support Vector Machine - with Soft Margin!

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers

$$0 \leq a_n \leq C$$

- if $a_n = 0$, then such points not used in prediction
- if $a_n < C$ because $a_n = C - \mu_n$ implies $\mu_n > 0$, from $\mu_n \xi_n = 0$ we get $\xi_n = 0$
- if $a_n = C$, then such points lie inside the margin



memory-based methods

Today

perhaps most powerful regressor
classification before deep learning

supervised
discriminative

perhaps most powerful classifier
regression before deep learning

Kernel Methods

Gaussian Process

Support Vector Machine

Dimension Reduction

actually Gaussian Process is not limited to regression
Support Vector Machine is not limited to classification

generative

Clustering

HMM

Graphical Models

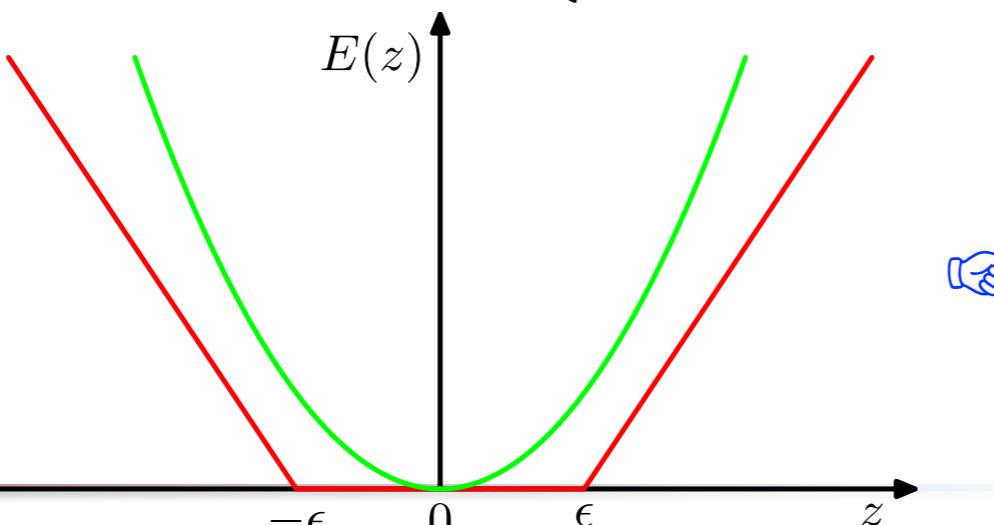
Support Vector Machine for Regression

- In simple linear regression, we minimise a regularised error function given by

$$\frac{1}{2} \sum_{n=1}^N \{y_n - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- To obtain sparse solution, we define an ϵ -insensitive error function which gives zero error if the absolute difference between the prediction $y(\mathbf{x})$ and the target t is less than ϵ where $\epsilon > 0$
 - Example:

$$E_\epsilon(y(\mathbf{x}) - t) = \begin{cases} 0, & \text{if } |y(\mathbf{x}) - t| < \epsilon; \\ |y(\mathbf{x}) - t| - \epsilon, & \text{otherwise} \end{cases}$$

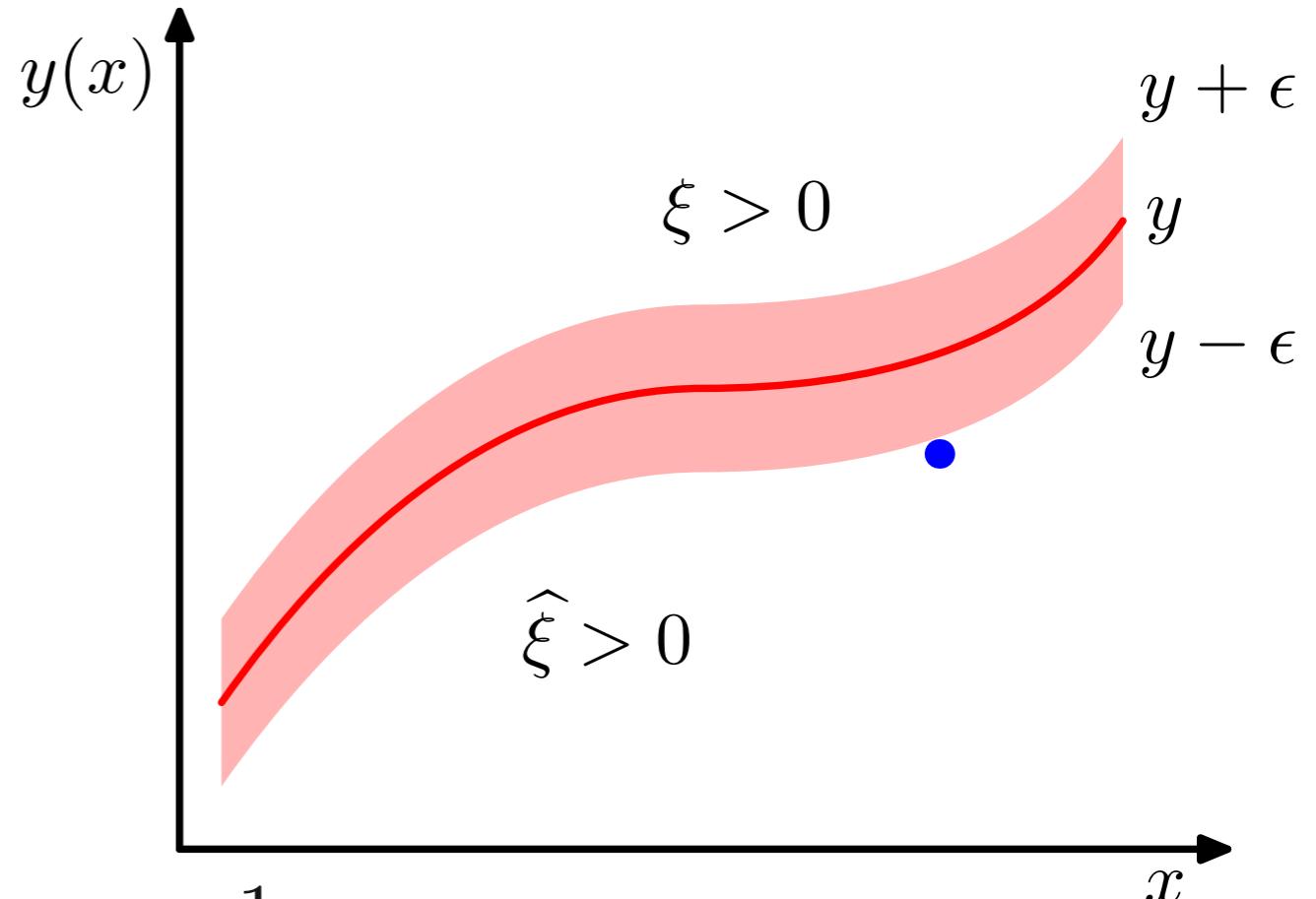


☞ $C \sum_{n=1}^N E_\epsilon(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2$

Support Vector Machine for Regression

- With also slack :)

Illustration of SVM regression, showing the regression curve together with the ϵ -insensitive ‘tube’. Also shown are examples of the slack variables ξ and $\hat{\xi}$. Points above the ϵ -tube have $\xi > 0$ and $\hat{\xi} = 0$, points below the ϵ -tube have $\xi = 0$ and $\hat{\xi} > 0$, and points inside the ϵ -tube have $\xi = \hat{\xi} = 0$.



$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } \xi_n \geq 0 \quad \hat{\xi}_n \geq 0$$

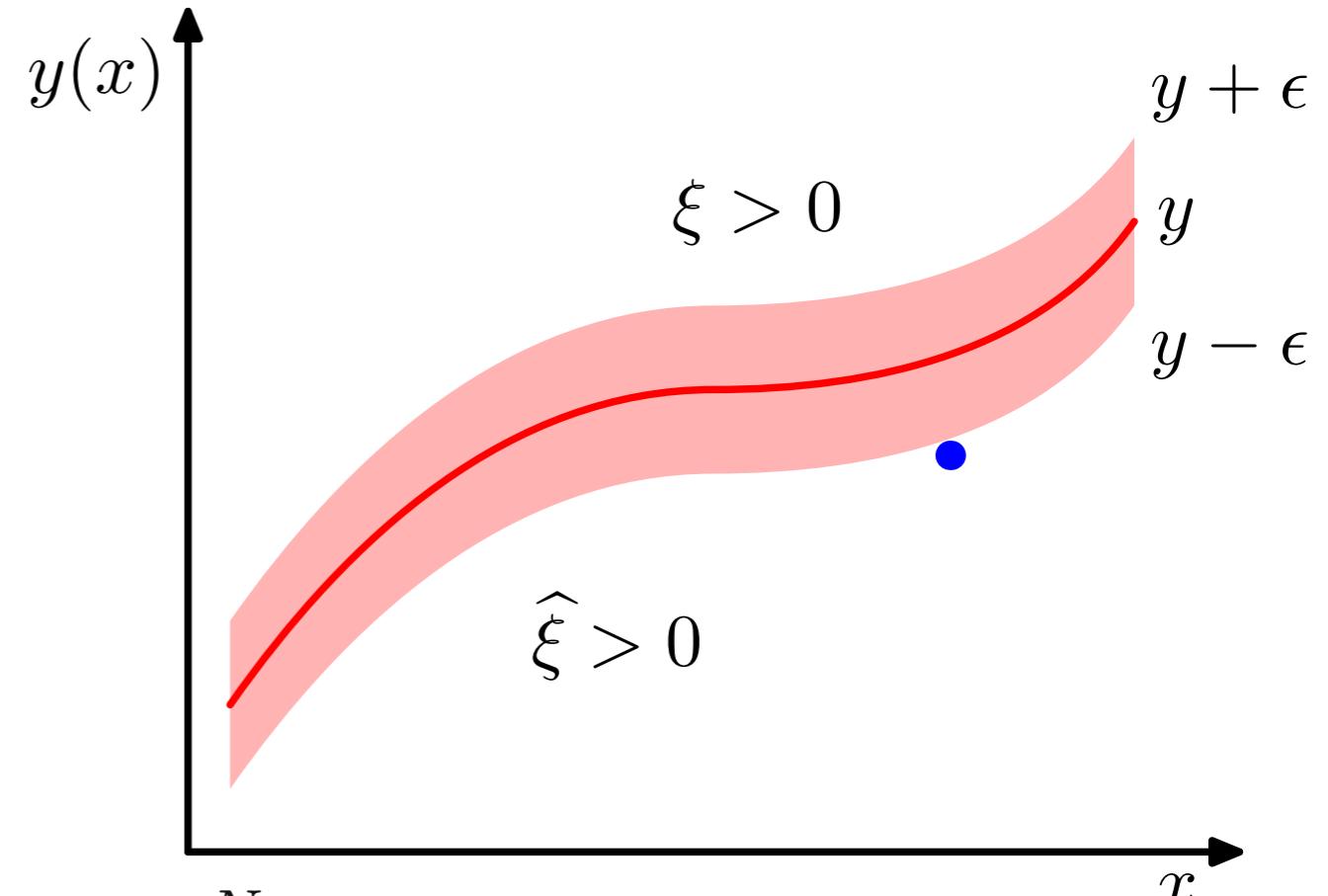
$$t_n \leq y(\mathbf{x}_n) + \epsilon + \xi_n$$

$$t_n \geq y(\mathbf{x}_n) - \epsilon - \hat{\xi}_n$$

Support Vector Machine for Regression

- With also slack :)

Illustration of SVM regression, showing the regression curve together with the ϵ -insensitive ‘tube’. Also shown are examples of the slack variables ξ and $\hat{\xi}$. Points above the ϵ -tube have $\xi > 0$ and $\hat{\xi} = 0$, points below the ϵ -tube have $\xi = 0$ and $\hat{\xi} > 0$, and points inside the ϵ -tube have $\xi = \hat{\xi} = 0$.



$$\begin{aligned}
 L = & C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) \\
 & - \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n - y_n + t_n)
 \end{aligned}$$

Support Vector Machine for Regression

- Dual representation
 - Maximising

$$\begin{aligned}\tilde{L}(\mathbf{a}, \hat{\mathbf{a}}) = & -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) \\ & - \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n\end{aligned}$$

subject to $0 \leq a_n \leq C$

$$0 \leq \hat{a}_n \leq C$$

$$\sum_{n=1}^N (a_n - \hat{a}_n) = 0$$

prediction: $y(\mathbf{x}) = \sum_{n=1}^N (a_n - \hat{a}_n) k(\mathbf{x}, \mathbf{x}_n) + b$

Support Vector Machine for Regression

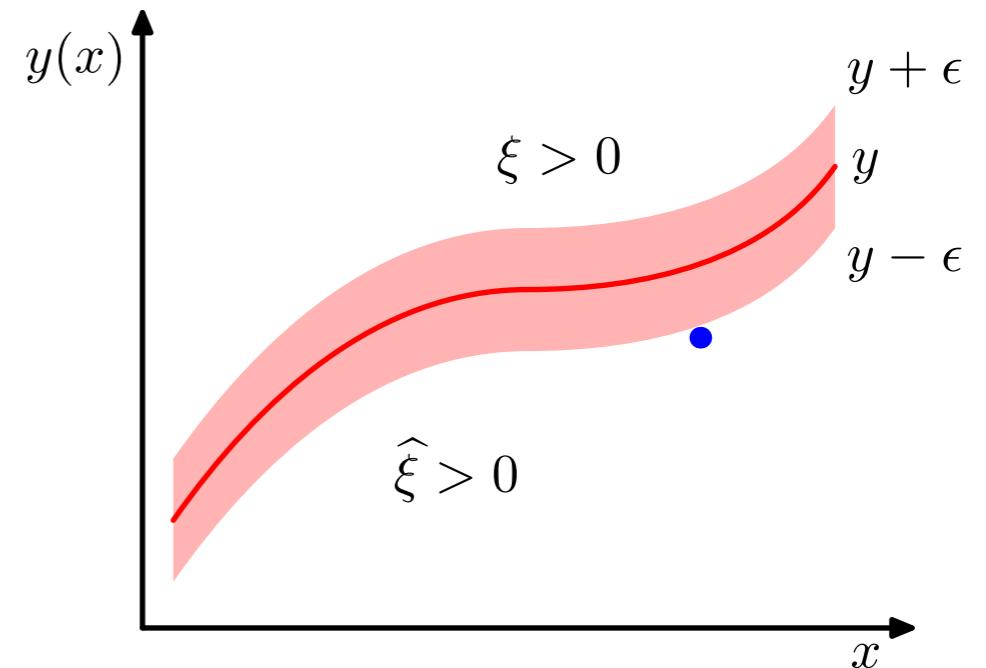
- KKT condition

$$a_n(\epsilon + \xi_n + y_n - t_n) = 0$$

$$\hat{a}_n(\epsilon + \hat{\xi}_n - y_n + t_n) = 0$$

$$(C - a_n)\xi_n = 0$$

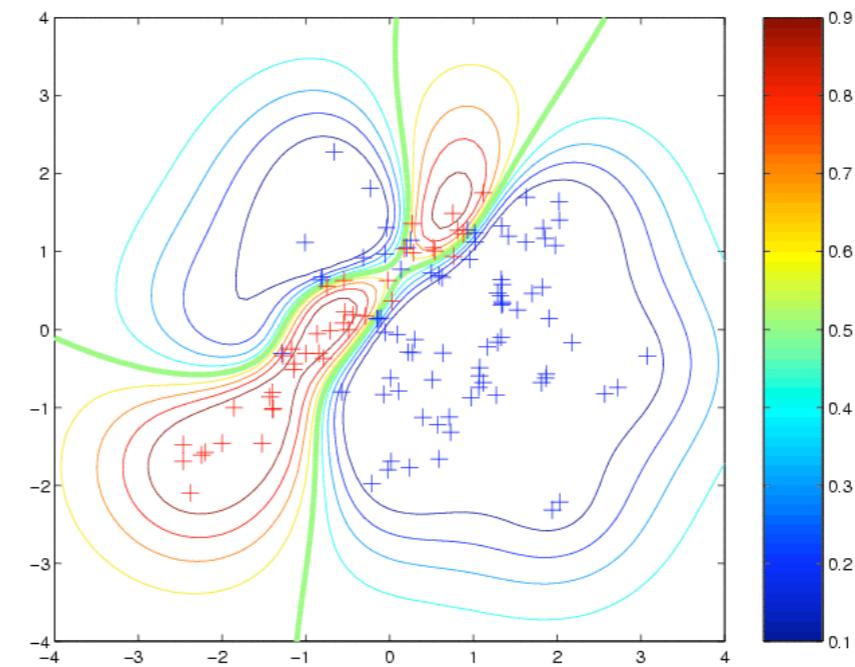
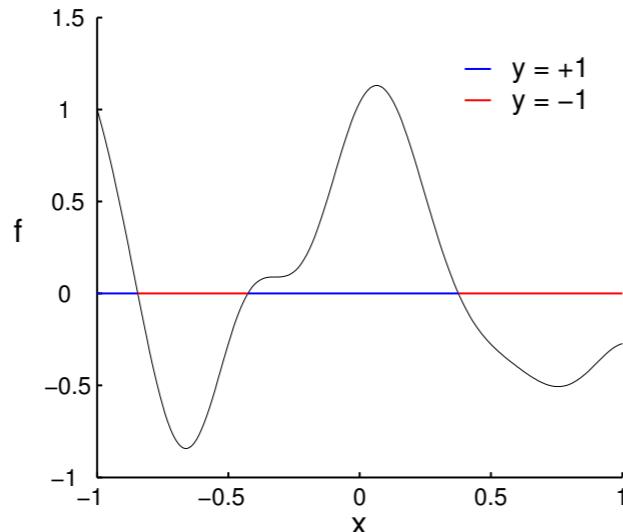
$$(C - \hat{a}_n)\hat{\xi}_n = 0.$$



- a_n can only be nonzero if $\epsilon + \xi_n + y_n - t_n = 0$
 - ☞ data points either lies on the upper boundary of the tube: $\xi_n = 0$, or above $\xi_n > 0$
- \hat{a}_n can only be nonzero if $\epsilon + \hat{\xi}_n - y_n + t_n = 0$
 - ☞ data points either lies on the lower boundary of the tube: $\hat{\xi}_n = 0$, or below $\hat{\xi}_n > 0$
- $\epsilon + \xi_n + y_n - t_n$ and $\epsilon + \hat{\xi}_n - y_n + t_n$ are incompatible
 - ☞ for every data point, either a_n or \hat{a}_n (or both) must be nonzero
- Support vectors are those data points either a_n nonzero or \hat{a}_n nonzero
 - ☞ these are data points that lie on the boundary of the tube or outside the tube
- All points within the tube have $a_n = \hat{a}_n = 0$

Gaussian Process for Classification

- Binary classification problem
 - Given dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{n=1}^n\} = (X, \mathbf{y})$ where \mathbf{y} are binary labels $\{-1, +1\}$, infer class label probabilities at new points.



- There are many ways to relate function values to class probabilities
- $$p(y_i|f_i) = \begin{cases} \frac{1}{1+\exp(-y_i f_i)} & \text{sigmoid (logistic)} \\ \Phi(y_i f_i) & \text{cumulative normal (probit)} \\ \mathbf{H}(y_i f_i) & \text{threshold} \\ \epsilon + (1 - 2\epsilon)\mathbf{H}(y_i f_i) & \text{robust threshold} \end{cases}$$
- Inference not that easy: approximation.

Limitations of SVM w.r.t. GPs

- The outputs of an SVM represent decisions rather than posterior probabilities.
- The SVM was originally formulated for two classes, and the extension to $K > 2$ classes is problematic.
- Gaussian Process can learn the kernel parameters automatically from data, no cross-validation is needed.
- Gaussian Process can be used for automatic feature selection.
- Gaussian Process can incorporate interpretable noise models and priors over functions, and can sample from prior to get intuitions about the model assumptions.