

Statistical (Machine) Learning

Lecture 11 Classifier Zoo
Part 2 of 2

Pierpaolo Brutti

Many routes to one destination...

...get close to that Bayes strategy!

- Empirical Risk Minimization vs Plug-in methods vs Local Averaging/Smoothing
- Parametric vs Semi-parametric vs Non-parametric
- Linear vs Non-linear
- Discriminative vs Generative
- Unconstrained vs Constrained vs Penalized/Regularized
- etc.

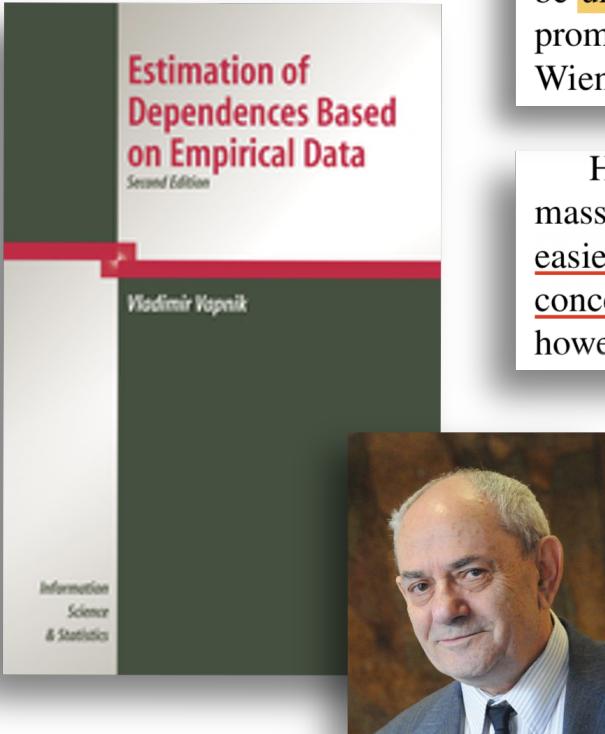
Many routes to one destination...

...get close to that Bayes strategy!

- Empirical Risk Minimization vs Plug-in methods vs Local Averaging/Smoothing
 - Parametric vs Semi-parametric vs Non-parametric
 - Linear vs Non-linear
 - Discriminative vs Generative
 - Unconstrained vs Constrained vs Penalized/Regularized
 - etc.
-

Focus more directly on:
Empirical Risk Minimization
(...although, clearly, also plug-in methods can be thought as ERM...)

In the beginning was...Vapnik

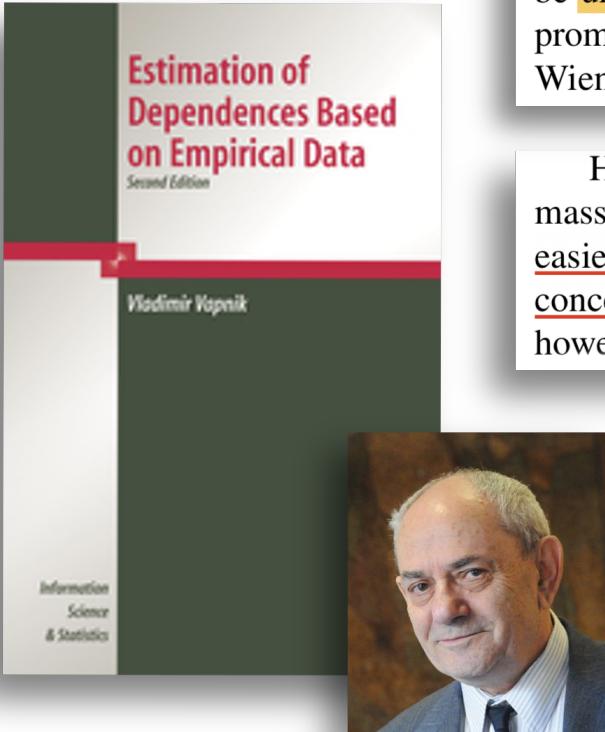


With the increasing role of science in everyday life, the general public began to discuss scientific discoveries in different areas: physical science, computer science (cybernetics), cognitive science (pattern recognition), biology (genomics), and philosophy. The discussions were held using very simplified scientific models that could be understood by the masses. Also scientists tried to appeal to the general public by promoting their philosophy using simplified models (for example, as has been done by Wiener). There is nothing wrong with this.

However, when science becomes a mass profession, the elements of the scientific mass culture in some cases start to substitute for the real scientific culture: It is much easier to learn the slogans of the scientific mass culture than it is to learn many different concepts from the original scientific sources. Science and “scientific mass culture,” however, are built on very different principles.

Gloss Over the Essentials and
Attract Attention to the Obvious

In the beginning was...Vapnik



With the increasing role of science in everyday life, the general public began to discuss scientific discoveries in different areas: physical science, computer science (cybernetics), cognitive science (pattern recognition), biology (genomics), and philosophy. The discussions were held using very simplified scientific models that could be understood by the masses. Also scientists tried to appeal to the general public by promoting their philosophy using simplified models (for example, as has been done by Wiener). There is nothing wrong with this.

However, when science becomes a mass profession, the elements of the scientific mass culture in some cases start to substitute for the real scientific culture: It is much easier to learn the slogans of the scientific mass culture than it is to learn many different concepts from the original scientific sources. Science and “scientific mass culture,” however, are built on very different principles.

Find the Essential
in the Nonobvious

Fast Forward: Toward Kernel-Machines

When solving a problem of interest, do not solve a more general problem as an intermediate step

- Logistic regression, being a *plug-in* method, constructs a (linear) classifier in two steps:
 1. First get an estimate the regression function $f_{\text{opt}}(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$
 2. Then place an example in class 1 if its estimated probability exceeds 0.5.

Fast Forward: Toward Kernel-Machines

When solving a problem of interest, do not solve a more general problem as an intermediate step

- Logistic regression, being a *plug-in* method, constructs a (linear) classifier in two steps:
 1. First get an estimate the regression function $f_{\text{opt}}(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$
 2. Then place an example in class 1 if its estimated probability exceeds 0.5.
- Gaussian Discriminant Analyses estimate the **entire** joint $p(y | \mathbf{x})$...even more extreme!

Fast Forward: Toward Kernel-Machines

When solving a problem of interest, do not solve a more general problem as an intermediate step

- Logistic regression, being a *plug-in* method, constructs a (linear) classifier in two steps:
 1. First get an estimate the regression function $f_{\text{opt}}(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$
 2. Then place an example in class 1 if its estimated probability exceeds 0.5.
- Gaussian Discriminant Analyses estimate the **entire** joint $p(y | \mathbf{x})$...even more extreme!
- In line with Vladimir Vapnik's "reductionism", SVM will completely avoids the first of the two steps and **geometrically** exploits the idea behind ERM with two crucial twists: **penalization** and **surrogate loss**

Fast Forward: Toward Kernel-Machines

When solving a problem of interest, do not solve a more general problem as an intermediate step

- Logistic regression, being a *plug-in* method, constructs a (linear) classifier in two steps:
 1. First get an estimate the regression function $f_{\text{opt}}(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$
 2. Then place an example in class 1 if its estimated probability exceeds 0.5.
- Gaussian Discriminant Analyses estimate the **entire** joint $p(y | \mathbf{x})$...even more extreme!
- In line with Vladimir Vapnik's "reductionism", SVM will completely avoids the first of the two steps and **geometrically** exploits the idea behind ERM with two crucial twists: **penalization** and **surrogate loss**
- In my opinion the *cool* aspect of this class of (linear) models is that all these *analytical* principles naturally flow from "simple" **geometric** arguments!

Fast Forward: Toward Kernel-Machines

When solving a problem of interest, do not solve a more general problem as an intermediate step

- Logistic regression, being a *plug-in* method, constructs a (linear) classifier in two steps:
 1. First get an estimate the regression function $f_{\text{opt}}(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$
 2. Then place an example in class 1 if its estimated probability exceeds 0.5.
- Gaussian Discriminant Analyses estimate the **entire** joint $p(y | \mathbf{x})$...even more extreme!
- In line with Vladimir Vapnik's "reductionism", SVM will completely avoids the first of the two steps and **geometrically** exploits the idea behind ERM with two crucial twists: **penalization** and **surrogate loss**
- In my opinion the *cool* aspect of this class of (linear) models is that all these *analytical* principles naturally flow from "simple" **geometric** arguments!
- ...then there is also the fact that everything can be made non-linear and extremely flexible using the so-called *kernel-trick*

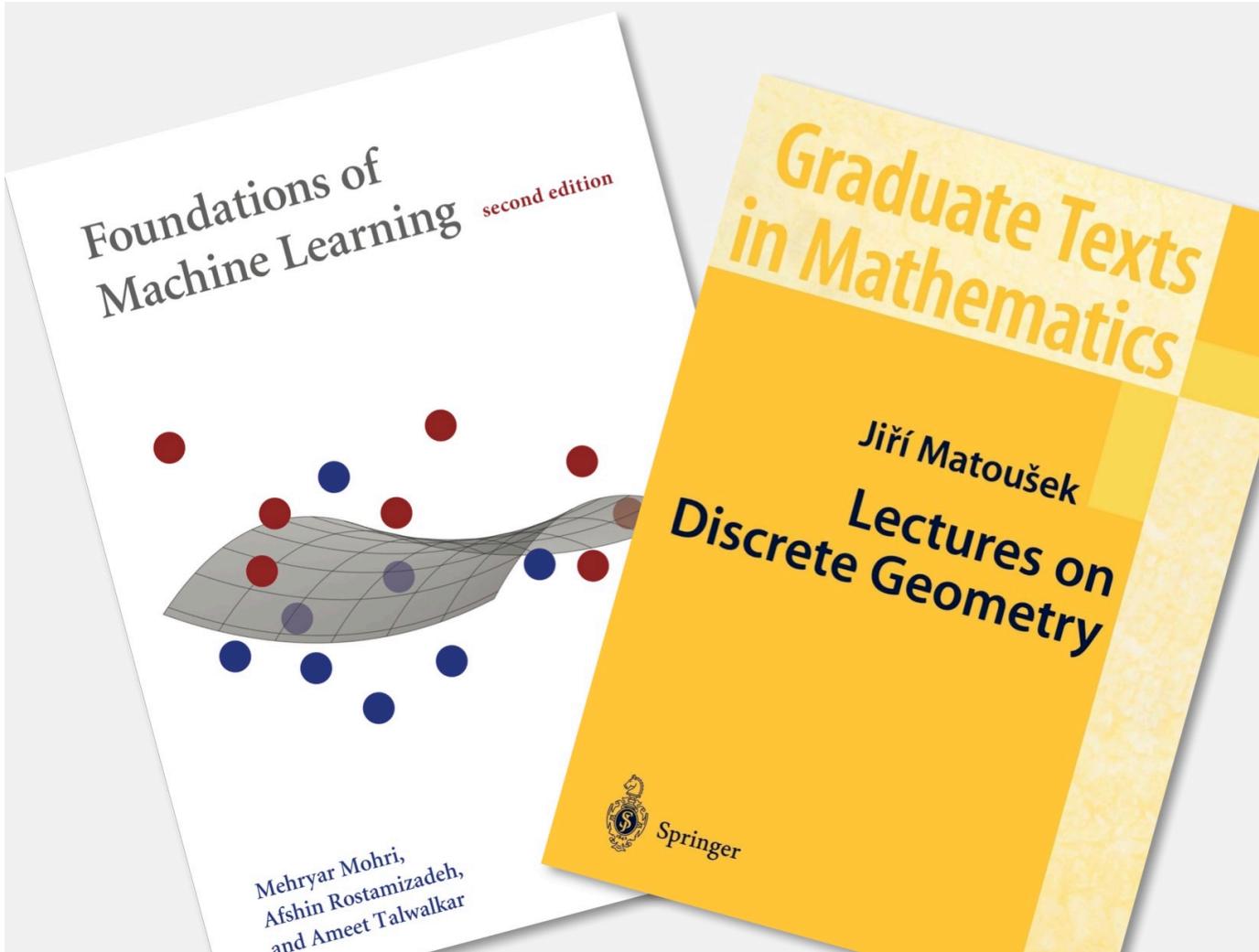
Fast Forward: Toward Kernel-Machines

When solving a problem of interest, do not solve a more general problem as an intermediate step

- Logistic regression, being a *plug-in* method, constructs a (linear) classifier in two steps:
 1. First get an estimate the regression function $f_{\text{opt}}(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$
 2. Then place an example in class 1 if its estimated probability exceeds 0.5.
- Gaussian Discriminant Analyses estimate the **entire** joint $p(y | \mathbf{x})$...even more extreme!
- In line with Vladimir Vapnik's "reductionism", SVM will completely avoids the first of the two steps and **geometrically** exploits the idea behind ERM with two crucial twists: **penalization** and **surrogate loss**
- In my opinion the *cool* aspect of this class of (linear) models is that all these *analytical* principles naturally flow from "simple" **geometric** arguments!
- ...then there is also the fact that everything can be made non-linear and extremely flexible using the so-called *kernel-trick*
- ..finally I should *also* mention that the whole construction comes with solid theoretical foundations: Vapnik–Chervonenkis theory (...zombie-Friday?!?)

...I mean...come on! Can you ask for more?

Fast Forward: Toward Kernel-Machines



Fast Forward: Toward Kernel-Machines

Going Non-Linear: Previously

- **Linear functions in some explicit transformed feature space**
 - Define **explicitly** a feature map $\phi : \mathcal{X} \mapsto \mathbb{R}^d$, we consider $f(\mathbf{x}) = \boldsymbol{\beta}^T \phi(\mathbf{x})$, with parameters $\boldsymbol{\beta} \in \mathbb{R}^d$ where d is the embedding dimension (possibly pretty large but to be controlled \rightsquigarrow overfit, think *polynomials*).
 - **Pros:** simple to implement, convex optimization with gradient descent algorithms, with running time complexity in $O(nd)$, and theoretical guarantees.
 - **Cons:** only applies to linear functions on **explicit** and **fixed** feature spaces, so they *can* underfit.

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Previously

- **Linear functions in some explicit transformed feature space**
 - Define **explicitly** a feature map $\phi : \mathcal{X} \mapsto \mathbb{R}^d$, we consider $f(\mathbf{x}) = \boldsymbol{\beta}^T \phi(\mathbf{x})$, with parameters $\boldsymbol{\beta} \in \mathbb{R}^d$ where d is the embedding dimension (possibly pretty large but to be controlled \rightsquigarrow overfit, think *polynomials*).
 - **Pros:** simple to implement, convex optimization with gradient descent algorithms, with running time complexity in $O(nd)$, and theoretical guarantees.
 - **Cons:** only applies to linear functions on **explicit** and **fixed** feature spaces, so they *can* underfit.
- **Neural Network**
 - Introduce non-linearity by chaining together simple layers that will implement a form of progressive non-linear data "distillation".
 - **Pros:** highly flexible and, at this point in time, surprisingly scalable given the available computing infrastructures!
 - **Cons:** these models are associated with non-convex optimization problems where gradient descent algorithms can be applied without guarantees, although several tricks have been observed to lead to better stability and performance.

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Now

- **Linear functions in some implicit feature space through kernel methods**
 - Define the map ϕ **implicitly** via a *kernel* function: a **similarity** measure between \mathbf{x} 's if you wish.
 - The feature map can have arbitrarily large dimension: ϕ belongs to an Hilbert space \mathcal{H} and we play with it "indirectly" through a *kernel* $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Now

- **Linear functions in some implicit feature space through kernel methods**
 - Define the map ϕ **implicitly** via a *kernel* function: a **similarity** measure between \mathbf{x} 's if you wish.
 - The feature map can have arbitrarily large dimension: ϕ belongs to an Hilbert space \mathcal{H} and we play with it "indirectly" through a *kernel* $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$
 - **Pros:** non-linear and flexible, easily adaptable to non-euclidean \mathcal{X} -space (just need a *similarity measure*) \rightsquigarrow be compared with NN-based embeddings (e.g. for text, **BERT** and **similia**, **FastText**, etc.)

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Now

- **Linear functions in some implicit feature space through kernel methods**
 - Define the map ϕ **implicitly** via a *kernel* function: a **similarity** measure between \mathbf{x} 's if you wish.
 - The feature map can have arbitrarily large dimension: ϕ belongs to an Hilbert space \mathcal{H} and we play with it "indirectly" through a *kernel* $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$
 - **Pros:** non-linear and flexible, easily adaptable to non-euclidean \mathcal{X} -space (just need a *similarity measure*) ↪ be compared with NN-based embeddings (e.g. for text, **BERT** and **similia**, **FastText**, etc.)
 - **Pros:** simple to implement, convex optimization algorithms with strong guarantees. Easily adapt to the regularity of the target function.

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Now

- **Linear functions in some implicit feature space through kernel methods**
 - Define the map ϕ **implicitly** via a *kernel* function: a **similarity** measure between \mathbf{x} 's if you wish.
 - The feature map can have arbitrarily large dimension: ϕ belongs to an Hilbert space \mathcal{H} and we play with it "indirectly" through a *kernel* $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$
 - **Pros:** non-linear and flexible, easily adaptable to non-euclidean \mathcal{X} -space (just need a *similarity measure*) \rightsquigarrow be compared with NN-based embeddings (e.g. for text, **BERT** and **similia**, **FastText**, etc.)
 - **Pros:** simple to implement, convex optimization algorithms with strong guarantees. Easily adapt to the regularity of the target function.
 - **Cons:** running-time complexity up to $O(n^2)$ (\rightsquigarrow *random sinks* \rightsquigarrow **Rahimi vs LeCun**) may still suffer from the curse of dimensionality for non-smooth target functions.

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Now

- **Linear functions in some implicit feature space through kernel methods**
 - Define the map ϕ **implicitly** via a *kernel* function: a **similarity** measure between \mathbf{x} 's if you wish.
 - The feature map can have arbitrarily large dimension: ϕ belongs to an Hilbert space \mathcal{H} and we play with it "indirectly" through a *kernel* $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$
 - **Pros:** non-linear and flexible, easily adaptable to non-euclidean \mathcal{X} -space (just need a *similarity measure*) \rightsquigarrow be compared with NN-based embeddings (e.g. for text, **BERT** and **similia**, **FastText**, etc.)
 - **Pros:** simple to implement, convex optimization algorithms with strong guarantees. Easily adapt to the regularity of the target function.
 - **Cons:** running-time complexity up to $O(n^2)$ (\rightsquigarrow *random sinks* \rightsquigarrow **Rahimi vs LeCun**) may still suffer from the curse of dimensionality for non-smooth target functions.
- Historically speaking, as NN started to gain some traction in the 1990s, a new approach to machine learning rose to fame and quickly sent NN back to oblivion: *kernel methods* wiped out NN!
- Kernel methods are a group of techniques, the best known of which is the support vector machine (SVM), whose modern formulation was developed by **Vladimir Vapnik** and **Corinna Cortes** in the early 1990s at Bell Labs and published in 1995, although an older linear formulation was published by Vapnik and **Alexey Chervonenkis** as early as 1963.

Fast Forward: Toward Kernel-Machines

Going Non-Linear: Now

- **Linear functions in some implicit feature space through kernel methods**
 - Define the map ϕ **implicitly** via a *kernel* function: a **similarity** measure between \mathbf{x} 's if you wish.
 - The feature map can have arbitrarily large dimension: ϕ belongs to an Hilbert space \mathcal{H} and we play with it "indirectly" through a *kernel* $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$
 - **Pros:** non-linear and flexible, easily adaptable to non-euclidean \mathcal{X} -space (just need a *similarity measure*) \rightsquigarrow be compared with NN-based embeddings (e.g. for text, **BERT** and **similia**, **FastText**, etc.)
 - **Pros:** simple to implement, convex optimization algorithms with strong guarantees. Easily adapt to the regularity of the target function.
 - **Cons:** running-time complexity up to $O(n^2)$ (\rightsquigarrow *random sinks* \rightsquigarrow **Rahimi vs LeCun**) may still suffer from the curse of dimensionality for non-smooth target functions.
- Historically speaking, as NN started to gain some traction in the 1990s, a new approach to machine learning rose to fame and quickly sent NN back to oblivion: *kernel methods* wiped out NN!
- **Progress** generally follows a *sigmoid curve*: it starts with a period of fast progress, which gradually stabilizes as researchers hit hard limitations, and then further improvements become incremental. *Deep learning* seems to be in the first half of that sigmoid now, whereas *kernel methods* are lying flat: we'll see...but there are **clear, formal connections** between (*shallow*) NN and kernel machines.

Empirical Risk Minimization: Part 1 of 2

Support Vector Machines

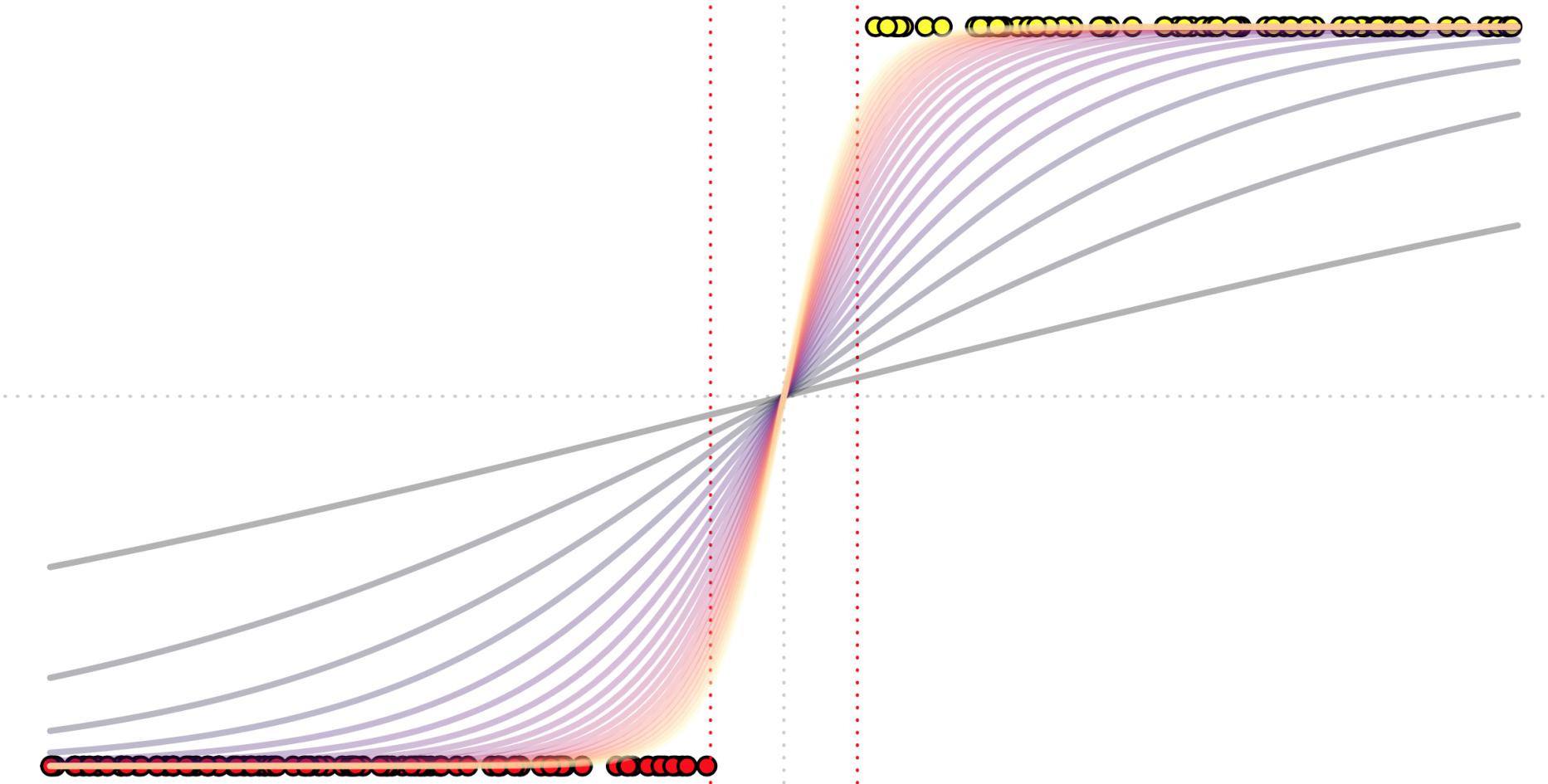
Level 1: Easy

The Linearly Separable Case

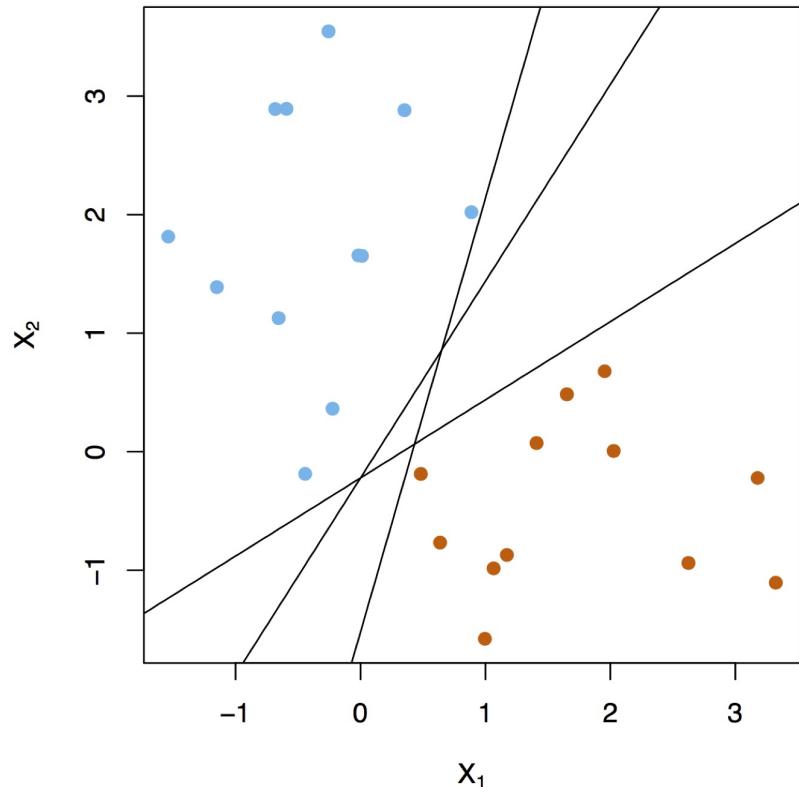
(Hard Margin)

Logistic Regression: The Dark Side

...let's start from an unexpected weakness of logistic regression...

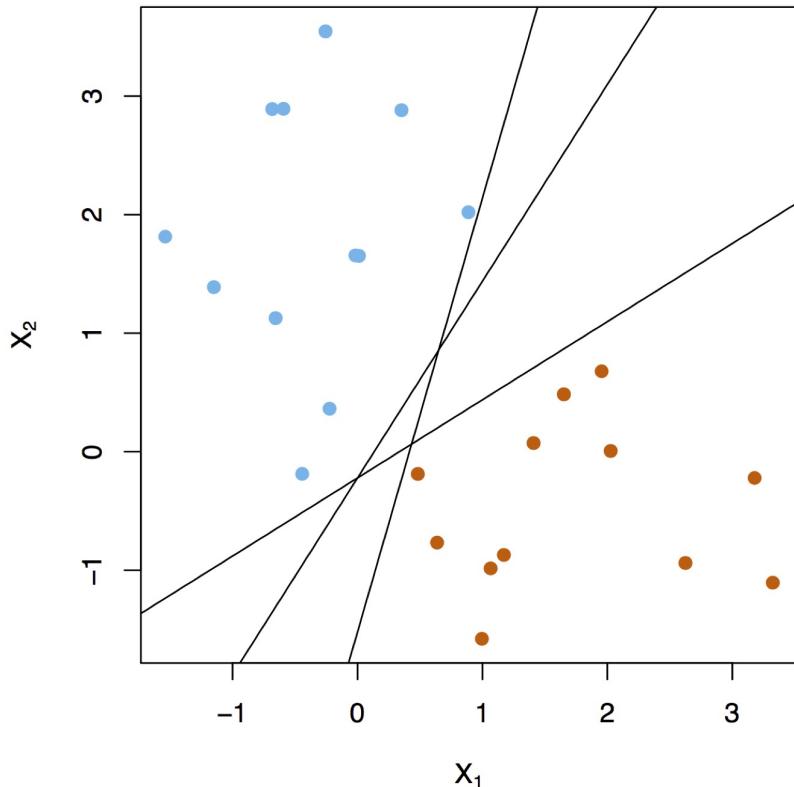


SVM: ...a question of margin!

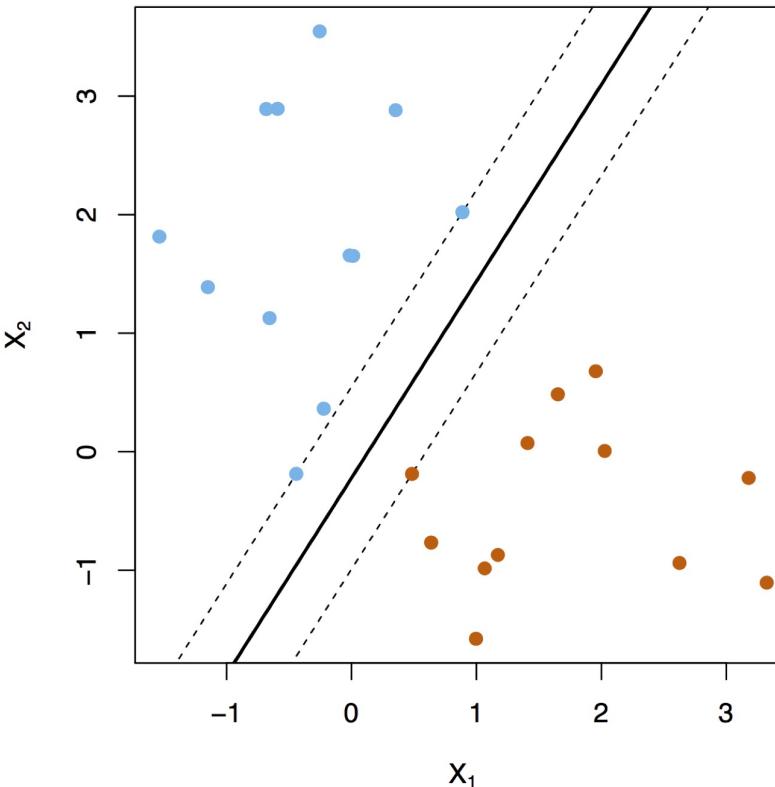


Filter out non-separating hyperplanes
(Hard-margin Constraint)

SVM: ...a question of margin!

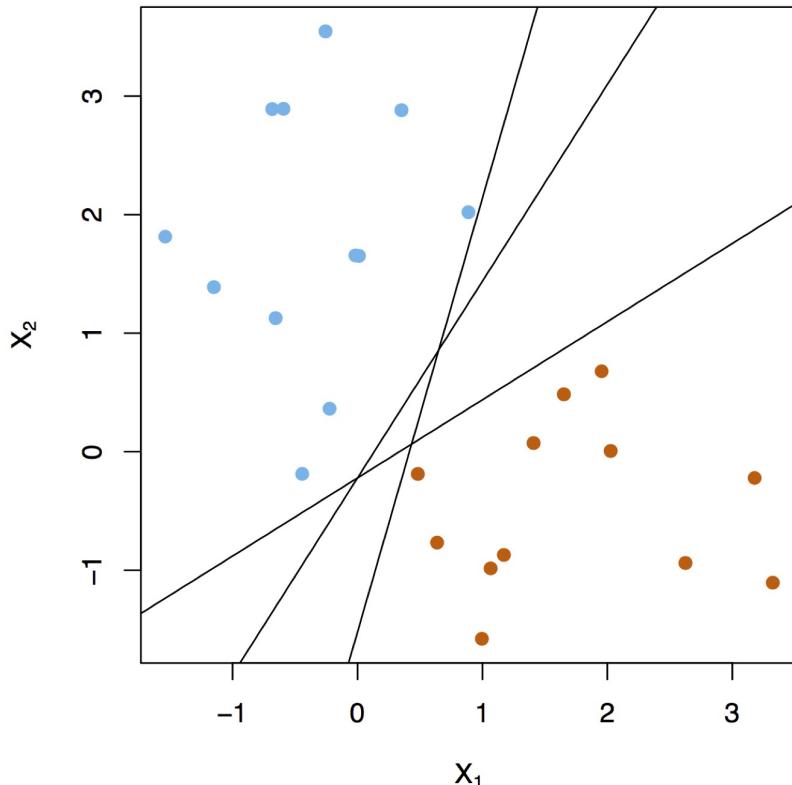


Filter out non-separating hyperplanes
(Hard-margin Constraint)

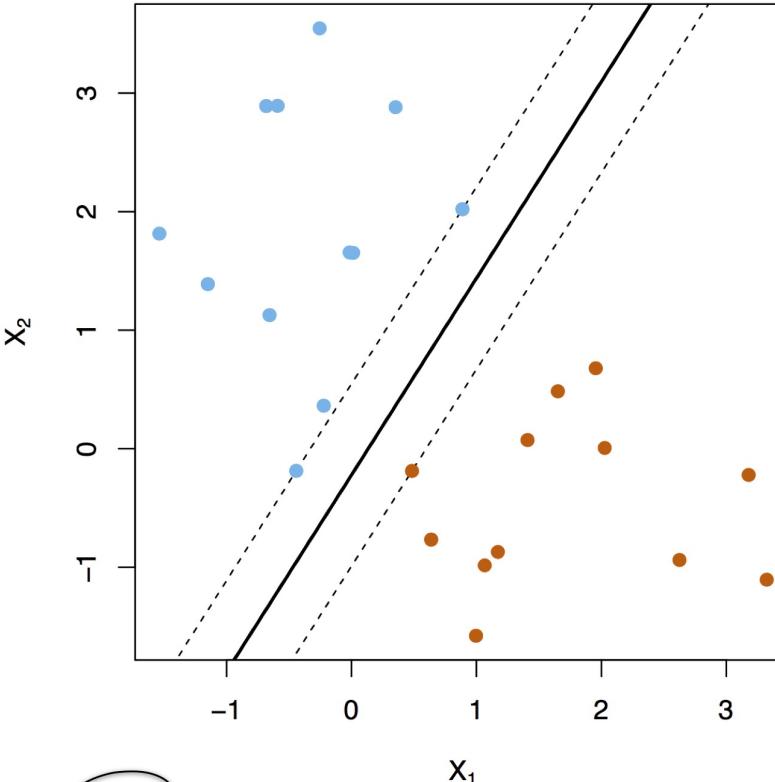


Pick the one with the largest “margin”
(Objective Function to Optimize)

SVM: ...a question of margin!



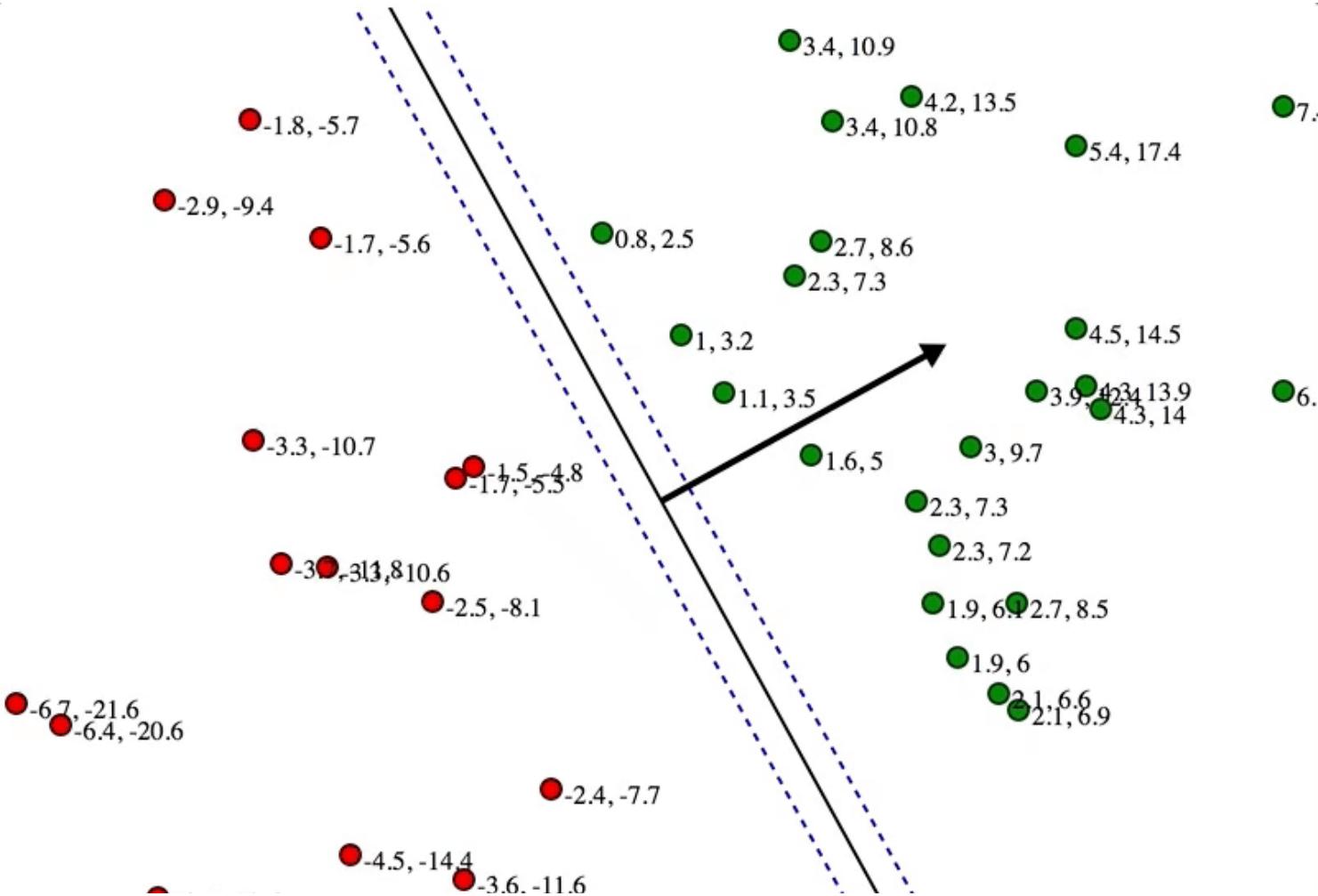
Filter out non-separating hyperplanes
(Hard-margin Constraint)



Pick the one with the largest “margin”
(Objective Function to Optimize)

Constrained (Convex) Optimization Problem

SVM: ...a question of margin!



Previously

- From now on, assume $Y \in \{-1, +1\}$, so a generic classifier can written as

$$h(\mathbf{x}) = \text{sign}(f(\mathbf{x})) \quad \text{with} \quad R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \quad \text{here} \quad f(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta} = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle_{\mathbb{R}^p}$$

Previously

- From now on, assume $Y \in \{-1, +1\}$, so a generic classifier can be written as

$$h(\mathbf{x}) = \text{sign}(f(\mathbf{x})) \quad \text{with} \quad R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \quad \text{here} \quad f(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta} = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle_{\mathbb{R}^p}$$

- Classifiers that compute a linear combination of the input features and return the sign, were called **perceptrons** in the engineering literature in the late 1950s ([Rosenblatt, 1958](#) on *Psychological Review*!).
- Perceptrons + their *online* training are another foundational brick in our modern learning construction

Psychological Review
Vol. 65, No. 6, 1958

THE PERCEPTRON: A PROBABILISTIC MODEL FOR INFORMATION STORAGE AND ORGANIZATION IN THE BRAIN¹

F. ROSENBLATT

Cornell Aeronautical Laboratory

If we are eventually to understand the capability of higher organisms for perceptual recognition, generalization, recall, and thinking, we must first have answers to three fundamental questions:

- How is information about the physical world sensed, or detected, by the biological system?
- In what form is information stored, or remembered?

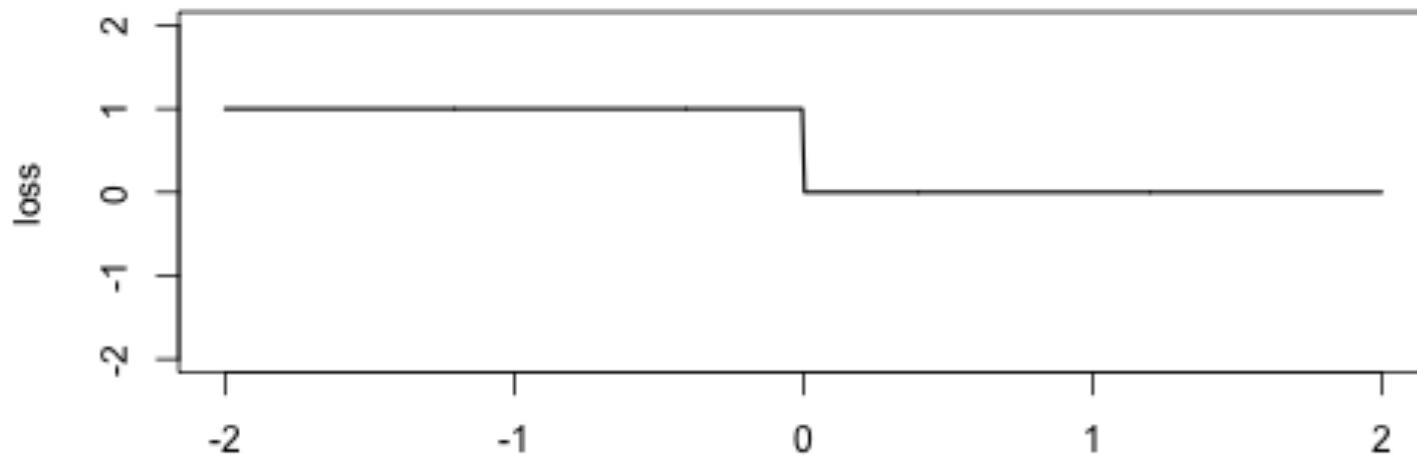
and the stored pattern. According to this hypothesis, if one understood the code or "wiring diagram" of the nervous system, one should, in principle, be able to discover exactly what an organism remembers by reconstructing the original sensory patterns from the "memory traces" which they have left, much as we might develop a photographic negative, or translate the pattern of electrical charges in the "memory" of a digital computer.

Previously

- From now on, assume $Y \in \{-1, +1\}$, so a generic classifier can be written as

$$h(\mathbf{x}) = \text{sign}(f(\mathbf{x})) \quad \text{with} \quad R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \quad \text{here} \quad f(\mathbf{x}) = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta} = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle_{\mathbb{R}^p}$$

- The *natural* loss is the 0-1 loss but this is annoying to optimize as not differentiable \rightsquigarrow *surrogates*



- Idea:** come up (geometrically!) with a smoother loss that *approximates/bounds* 0-1 loss \rightsquigarrow the *hinge*

SVM: ...a question of margin(s)!

Confidence or Functional Margin

- $R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \rightsquigarrow$ getting the class wrong is **not** so bad if we are *near* the decision boundary
 - ...getting the class right is *better* if we are **far** from the decision boundary

SVM: ...a question of margin(s)!

Confidence or Functional Margin

- $R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \rightsquigarrow$ getting the class wrong is **not** so bad if we are *near* the decision boundary
 - ...getting the class right is *better* if we are **far** from the decision boundary
- This notion of **margin** also known as **confidence or functional margin**, is $y \cdot f(\mathbf{x})$
 - > 0 if the point (y, \mathbf{x}) is correctly classified
 - < 0 if the point (y, \mathbf{x}) is incorrectly classified
 - Bigger (margin) is better

SVM: ...a question of margin(s)!

Confidence or Functional Margin

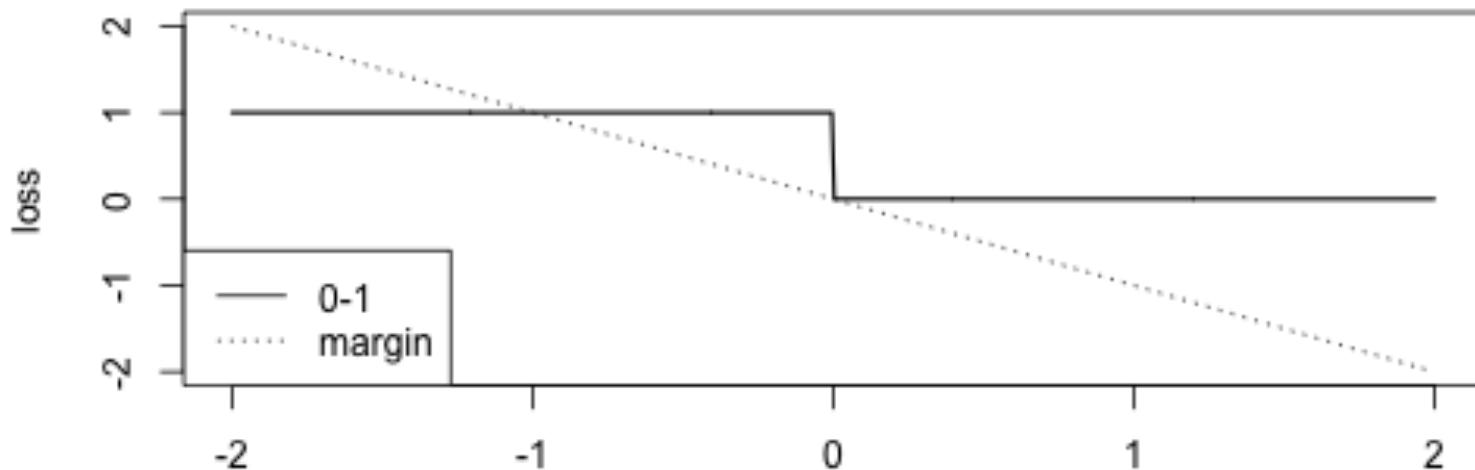
- $R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \rightsquigarrow$ getting the class wrong is **not** so bad if we are *near* the decision boundary
 - ...getting the class right is *better* if we are **far** from the decision boundary
- This notion of **margin** also known as **confidence or functional margin**, is $y \cdot f(\mathbf{x})$
 - > 0 if the point (y, \mathbf{x}) is correctly classified
 - < 0 if the point (y, \mathbf{x}) is incorrectly classified
 - Bigger (margin) is better
- The **margin loss** is $-y \cdot f(\mathbf{x})$
 - Keep with our minimize-all-the-things habits

SVM: ...a question of margin(s)!

Confidence or Functional Margin

- $R(h) = \mathbb{E} \Phi_{0-1}(Y \cdot f(\mathbf{X})) \rightsquigarrow$ getting the class wrong is **not** so bad if we are *near* the decision boundary
 - ...getting the class right is *better* if we are **far** from the decision boundary
- This notion of **margin** also known as **confidence or functional margin**, is $y \cdot f(\mathbf{x})$
 - > 0 if the point (y, \mathbf{x}) is correctly classified
 - < 0 if the point (y, \mathbf{x}) is incorrectly classified
 - Bigger (margin) is better
- The **margin loss** is $-y \cdot f(\mathbf{x})$
 - Keep with our minimize-all-the-things habits
- Margin at data point i is: $\rho_i = y_i \cdot f(\mathbf{x}_i)$
- Over-all margin \rightsquigarrow driven by the data-point closest to the boundary: $M(f) = \min_i \rho_i$

0-1 loss vs margin loss



SVM: ...a question of margin(s)!

Geometric Margin

- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(cf(\mathbf{x}))$ for any $c > 0$

SVM: ...a question of margin(s)!

Geometric Margin

- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(c f(\mathbf{x}))$ for any $c > 0$
- **Implication:** the weights $\boldsymbol{\beta}$ and $c \boldsymbol{\beta}$ give us the same (linear) classifier

SVM: ...a question of margin(s)!

Geometric Margin

- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(c f(\mathbf{x}))$ for any $c > 0$
- **Implication:** the weights $\boldsymbol{\beta}$ and $c \boldsymbol{\beta}$ give us the same (linear) classifier
- But the (confidence) margin, as defined, goes up if $c > 1$
 - Why not find a good plane and then "scale up" its $\boldsymbol{\beta}$ as far as the computer allows? \rightsquigarrow *non-sense*

SVM: ...a question of margin(s)!

Geometric Margin

- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(c f(\mathbf{x}))$ for any $c > 0$
- **Implication:** the weights $\boldsymbol{\beta}$ and $c \boldsymbol{\beta}$ give us the same (linear) classifier
- But the (confidence) margin, as defined, goes up if $c > 1$
 - Why not find a good plane and then "scale up" its $\boldsymbol{\beta}$ as far as the computer allows? \rightsquigarrow *non-sense*
- Two responses: **geometric margin**, or **fixed parameter scaling**

SVM: ...a question of margin(s)!

Geometric Margin

- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(c f(\mathbf{x}))$ for any $c > 0$
- **Implication:** the weights $\boldsymbol{\beta}$ and $c \boldsymbol{\beta}$ give us the same (linear) classifier
- But the (confidence) margin, as defined, goes up if $c > 1$
 - Why not find a good plane and then "scale up" its $\boldsymbol{\beta}$ as far as the computer allows? \rightsquigarrow *non-sense*
- Two responses: **geometric margin**, or **fixed parameter scaling**
- **Geometric margin** (see page 150 of [ESL](#))
Like margin, but *geometrically* equivalent to the (signed) **distance** from the (linear) boundary:

$$y \cdot \left(\frac{\beta_0}{\|\boldsymbol{\beta}\|_2} + \left\langle \frac{\boldsymbol{\beta}}{\|\boldsymbol{\beta}\|_2}, \mathbf{x} \right\rangle \right) = \frac{y \cdot (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle)}{\|\boldsymbol{\beta}\|_2}$$

SVM: ...a question of margin(s)!

Geometric Margin

- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(c f(\mathbf{x}))$ for any $c > 0$
- **Implication:** the weights $\boldsymbol{\beta}$ and $c \boldsymbol{\beta}$ give us the same (linear) classifier
- But the (confidence) margin, as defined, goes up if $c > 1$
 - Why not find a good plane and then "scale up" its $\boldsymbol{\beta}$ as far as the computer allows? \rightsquigarrow *non-sense*
- Two responses: **geometric margin**, or **fixed parameter scaling**
- **Geometric margin** (see page 150 of [ESL](#))
Like margin, but *geometrically* equivalent to the (signed) **distance** from the (linear) boundary:

$$y \cdot \left(\frac{\beta_0}{\|\boldsymbol{\beta}\|_2} + \left\langle \frac{\boldsymbol{\beta}}{\|\boldsymbol{\beta}\|_2}, \mathbf{x} \right\rangle \right) = \frac{y \cdot (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle)}{\|\boldsymbol{\beta}\|_2}$$

- **Fixed parameter scaling:** insist that $\|\boldsymbol{\beta}\|_2 = 1$
- **Notice (A):** when actually $\|\boldsymbol{\beta}\|_2 = 1 \rightsquigarrow$ confidence margin = geometric margin

SVM: ...a question of margin(s)!

Geometric Margin

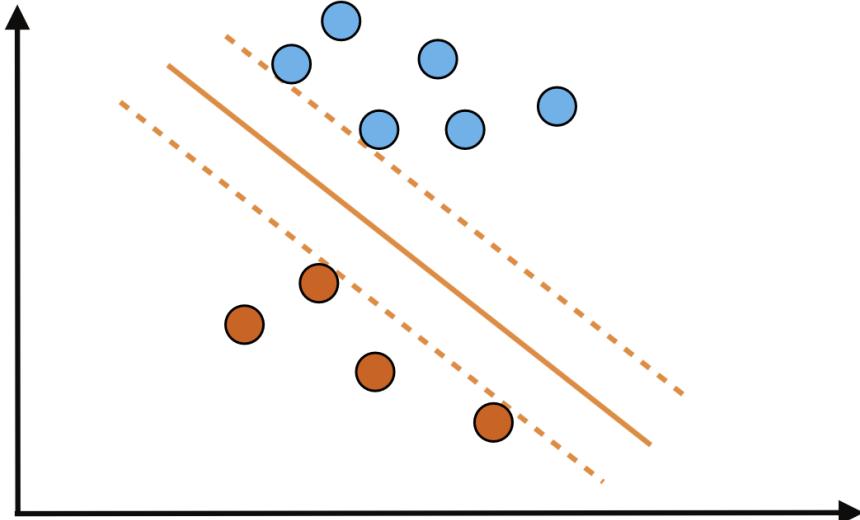
- **Notice:** $\text{sign}(f(\mathbf{x})) = \text{sign}(c f(\mathbf{x}))$ for any $c > 0$
- **Implication:** the weights β and $c\beta$ give us the same (linear) classifier
- But the (confidence) margin, as defined, goes up if $c > 1$
 - Why not find a good plane and then "scale up" its β as far as the computer allows? \rightsquigarrow *non-sense*
- Two responses: **geometric margin**, or **fixed parameter scaling**
- **Geometric margin** (see page 150 of [ESL](#))
Like margin, but *geometrically* equivalent to the (signed) **distance** from the (linear) boundary:

$$y \cdot \left(\frac{\beta_0}{\|\beta\|_2} + \left\langle \frac{\beta}{\|\beta\|_2}, \mathbf{x} \right\rangle \right) = \frac{y \cdot (\beta_0 + \langle \beta, \mathbf{x} \rangle)}{\|\beta\|_2}$$

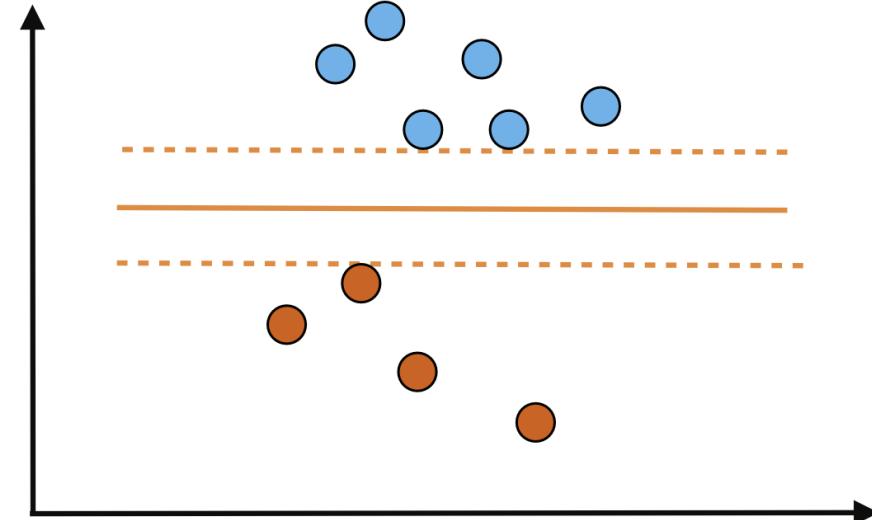
- **Fixed parameter scaling:** insist that $\|\beta\|_2 = 1$
- **Notice (B):** different scaling/norms \rightsquigarrow different margin \rightsquigarrow different algorithm!

SVM: ...a question of margin(s)!

(Other) Geometric Margin(s)



Norm $\|\cdot\|_2$
(Linear SVM)



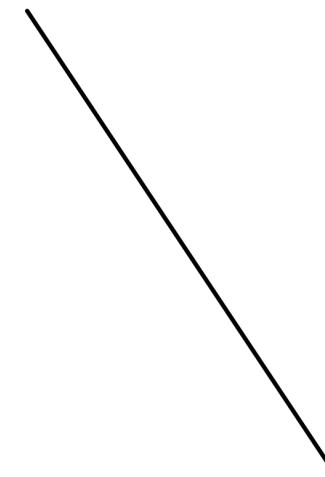
Norm $\|\cdot\|_1$
(Linear Ensembles \rightsquigarrow AdaBoost)

SVM: ...a question of margin(s)!

Geometric Margin

Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes

$$\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle = 0$$

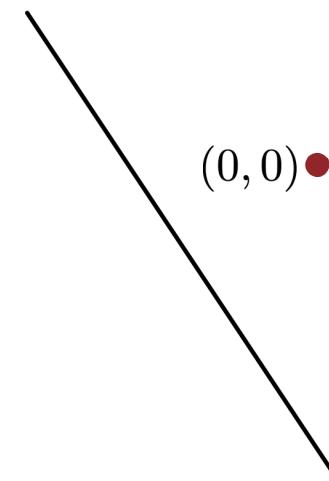


SVM: ...a question of margin(s)!

Geometric Margin

Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes

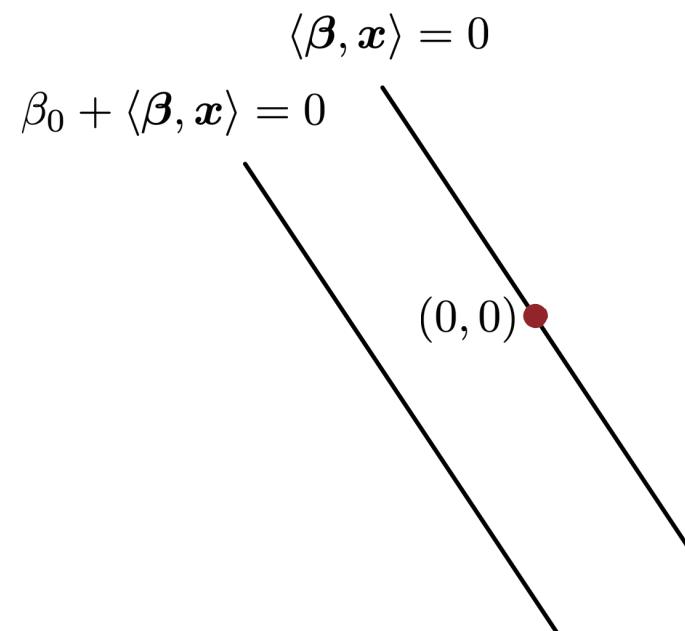
$$\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle = 0$$



SVM: ...a question of margin(s)!

Geometric Margin

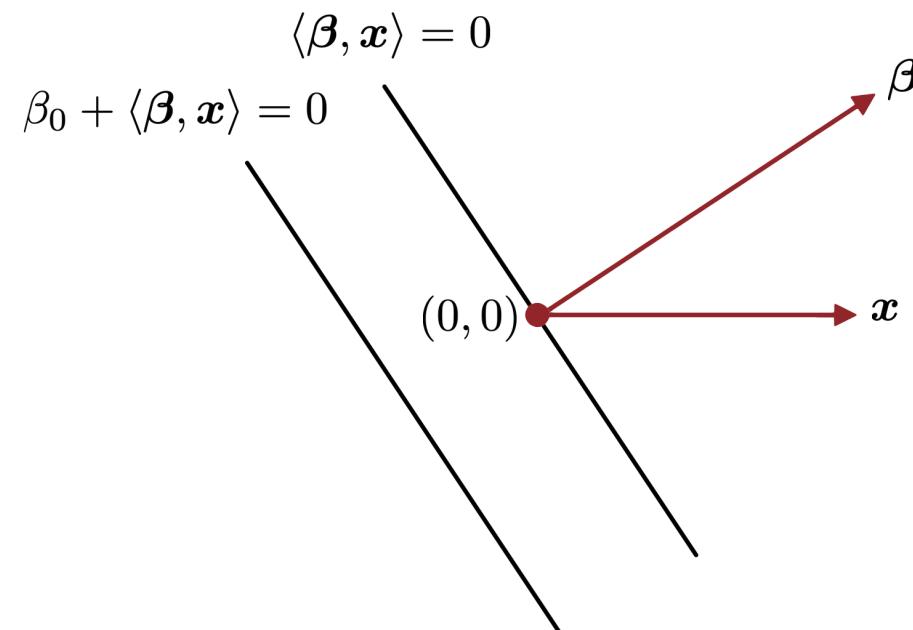
Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes



SVM: ...a question of margin(s)!

Geometric Margin

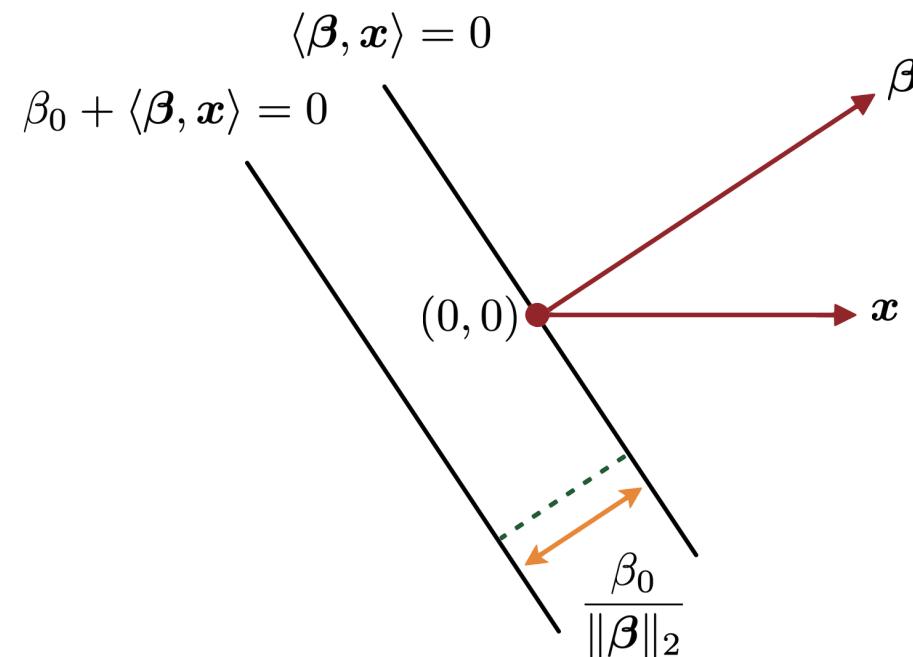
Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes



SVM: ...a question of margin(s)!

Geometric Margin

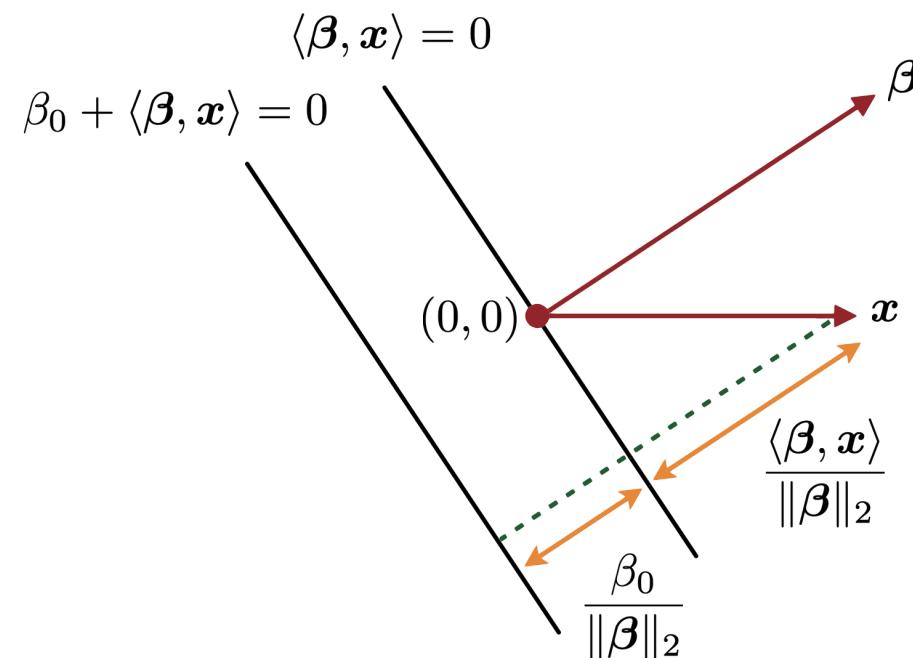
Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes



SVM: ...a question of margin(s)!

Geometric Margin

Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes



SVM: ...a question of margin(s)!

Geometric Margin

Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes

- The **geometric margin** $\rho(\mathbf{x})$ at a point \mathbf{x} is its **Euclidean distance** to the hyperplane $\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle = 0$

$$\rho(\mathbf{x}) = \frac{|\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle|}{\|\boldsymbol{\beta}\|_2}$$

SVM: ...a question of margin(s)!

Geometric Margin

Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes

- The **geometric margin** $\rho(\mathbf{x})$ at a point \mathbf{x} is its **Euclidean distance** to the hyperplane $\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle = 0$

$$\rho(\mathbf{x}) = \frac{|\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle|}{\|\boldsymbol{\beta}\|_2}$$

- The **geometric margin** ρ of a linear classifier for a (training) sample $\mathcal{D}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is the *minimum geometric margin* over the points in the sample

$$\rho = \min_i \rho(\mathbf{x}_i)$$

that is, the distance of the hyperplane defining the classifier to the closest sample points

SVM: ...a question of margin(s)!

Geometric Margin

Step into this other type of margin geometrically:
Signed **distance** from the (linear) boundary between classes

- The **geometric margin** $\rho(\mathbf{x})$ at a point \mathbf{x} is its **Euclidean distance** to the hyperplane $\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle = 0$

$$\rho(\mathbf{x}) = \frac{|\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle|}{\|\boldsymbol{\beta}\|_2}$$

- The **geometric margin** ρ of a linear classifier for a (training) sample $\mathcal{D}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is the *minimum geometric margin* over the points in the sample

$$\rho = \min_i \rho(\mathbf{x}_i)$$

that is, the distance of the hyperplane defining the classifier to the closest sample points

- The SVM solution is the separating hyperplane with the **maximum** (geometric) margin
 - It can be viewed as the **safest** choice: a **test** point is classified correctly by such a plane even when it falls within a distance ρ of the *training samples* sharing the same label.

SVM: ...a question of margin(s)!

Geometric Margin

Now a bit of algebra and thinking...

$$\rho = \max_{\substack{\text{All} \\ \text{Separating} \\ \text{Hyperplane}}} \min_i \rho(\mathbf{x}_i) = \max_{\substack{(\beta_0, \boldsymbol{\beta}) \\ y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 0}} \min_i \frac{|\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle|}{\|\boldsymbol{\beta}\|_2} = \max_{(\beta_0, \boldsymbol{\beta})} \frac{\min_i y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle)}{\|\boldsymbol{\beta}\|_2}$$

SVM: ...a question of margin(s)!

Geometric Margin

Now a bit of algebra and thinking...

$$\rho = \max_{\substack{\text{All} \\ \text{Separating} \\ \text{Hyperplane}}} \min_i \rho(\mathbf{x}_i) = \max_{\substack{(\beta_0, \boldsymbol{\beta}) \\ y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 0}} \min_i \frac{|\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle|}{\|\boldsymbol{\beta}\|_2} = \max_{(\beta_0, \boldsymbol{\beta})} \frac{\min_i y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle)}{\|\boldsymbol{\beta}\|_2}$$

- The 3rd equality follows from the fact that, since the sample is linearly separable, at the maximimizing pair $(\beta_0, \boldsymbol{\beta})$, $y_i \cdot (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle)$ must be non negative for each i .

SVM: ...a question of margin(s)!

Geometric Margin

Now a bit of algebra and thinking...

$$\rho = \max_{\substack{\text{All} \\ \text{Separating} \\ \text{Hyperplane}}} \min_i \rho(\mathbf{x}_i) = \max_{\substack{(\beta_0, \beta) \\ y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 0}} \min_i \frac{|\beta_0 + \langle \beta, \mathbf{x}_i \rangle|}{\|\beta\|_2} = \max_{(\beta_0, \beta)} \frac{\min_i y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle)}{\|\beta\|_2}$$

- The 3rd equality follows from the fact that, since the sample is linearly separable, at the maximimizing pair (β_0, β) , $y_i \cdot (\beta_0 + \langle \beta, \mathbf{x}_i \rangle)$ must be non negative for each i .
- **Observe:** the last expression is invariant to multiplication of (β_0, β) by any positive scalar c
- **Consequence:** Restrict the search to (β_0, β) scaled such that

$$\min_i y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle) = 1 \Leftrightarrow \min_i |\beta_0 + \langle \beta, \mathbf{x}_i \rangle| = 1$$

SVM: ...a question of margin(s)!

Geometric Margin

Now a bit of algebra and thinking...

$$\rho = \max_{\substack{\text{All} \\ \text{Separating} \\ \text{Hyperplane}}} \min_i \rho(\mathbf{x}_i) = \max_{\substack{(\beta_0, \beta) \\ y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 0}} \min_i \frac{|\beta_0 + \langle \beta, \mathbf{x}_i \rangle|}{\|\beta\|_2} = \max_{(\beta_0, \beta)} \frac{\min_i y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle)}{\|\beta\|_2}$$

- The 3rd equality follows from the fact that, since the sample is linearly separable, at the maximimizing pair (β_0, β) , $y_i \cdot (\beta_0 + \langle \beta, \mathbf{x}_i \rangle)$ must be non negative for each i .
- **Observe:** the last expression is invariant to multiplication of (β_0, β) by any positive scalar c
- **Consequence:** Restrict the search to (β_0, β) scaled such that

$$\min_i y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle) = 1 \Leftrightarrow \min_i |\beta_0 + \langle \beta, \mathbf{x}_i \rangle| = 1$$

$$\Rightarrow \rho = \max_{\substack{(\beta_0, \beta) \\ \min_i y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle) = 1}} \frac{1}{\|\beta\|_2} = \max_{\substack{(\beta_0, \beta) \\ \forall i : y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1}} \frac{1}{\|\beta\|_2} = \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$$

Learning classifiers by maximizing the margin

- Only consider "linear machines" $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$ with suitably scaled coefficients such that:
- **Maximize** the margin $M(f)$
- **Subject to:** separability condition

Learning classifiers by maximizing the margin

- Only consider "linear machines" $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$ with suitably scaled coefficients such that:
- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
- **Subject to:** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$

Learning classifiers by maximizing the margin

- Only consider "linear machines" $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$ with suitably scaled coefficients such that:
 - **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
 - **Subject to:** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$
-

Quick Convex Optimization Interlude

- Consider the constrained **convex** optimization problem

$$\min_{\theta} M(\theta) \quad \text{subject to} \quad \Omega_j(\theta) \leq 0 \quad j \in \{1, \dots, n\}$$

- Define the **Lagrangian**: $\mathcal{L}(\theta, \boldsymbol{\alpha}) = M(\theta) + \sum_j \alpha_j \Omega_j(\theta)$
- The **dual function** is defined by $F(\boldsymbol{\alpha}) = \inf_{\theta} \mathcal{L}(\theta, \boldsymbol{\alpha})$
- A central result in **convex optimization** is that (under some conditions verified in our application), the original problem can be solved by **maximizing** the dual F subject to:
$$\alpha_j \geq 0 \quad \text{for each } j \in \{1, \dots, n\}$$

Interlude: More on Lagrange Multipliers

Constrained optimization in general

- Get a bit more abstract for a moment and think about constrained optimization in general

$$\theta^* = \min_{\theta \in \Theta} M(\theta)$$

subject to

$$\Omega(\theta) \leq c$$

- How might we solve this?

Constrained optimization in general

- Get a bit more abstract for a moment and think about constrained optimization in general

$$\theta^* = \min_{\theta \in \Theta} M(\theta)$$

subject to

$$\Omega(\theta) \leq c$$

- How might we solve this?

1. Use the **constraint equation** $\Omega(\theta) = c$ to **eliminate a degree of freedom**

- In other words, write one coordinate in θ as a function of the others and of c
- Do **unconstrained** optimization over the remaining degrees of freedom

Constrained optimization in general

- Get a bit more abstract for a moment and think about constrained optimization in general

$$\theta^* = \min_{\theta \in \Theta} M(\theta)$$

subject to

$$\Omega(\theta) \leq c$$

- How might we solve this?

1. Use the **constraint equation** $\Omega(\theta) = c$ to **eliminate a degree of freedom**

- In other words, write one coordinate in θ as a function of the others and of c
- Do **unconstrained** optimization over the remaining degrees of freedom

2. Ok, but what about the \leq case?

Constrained optimization in general

- Get a bit more abstract for a moment and think about constrained optimization in general

$$\theta^* = \min_{\theta \in \Theta} M(\theta)$$

subject to

$$\Omega(\theta) \leq c$$

- How might we solve this?

1. Use the **constraint equation** $\Omega(\theta) = c$ to **eliminate a degree of freedom**

- In other words, write one coordinate in θ as a function of the others and of c
- Do **unconstrained** optimization over the remaining degrees of freedom

2. Ok, but what about the \leq case?

3. Add a new variable and do **unconstrained** optimization over a **larger** problem

Lagrange multipliers: Constraints \longleftrightarrow Penalties

- If we have one *equality* constraint, say $\Omega(\theta) = c$, we'd add a **Lagrange multiplier**:

$$\min_{\theta \in \Theta, \alpha \in \mathbb{R}} M(\theta) + \alpha (\Omega(\theta) - c) = \min_{\theta \in \Theta, \alpha \in \mathbb{R}} \mathcal{L}(\theta, \alpha)$$

Lagrange multipliers: Constraints \longleftrightarrow Penalties

- If we have one *equality* constraint, say $\Omega(\theta) = c$, we'd add a **Lagrange multiplier**:

$$\min_{\theta \in \Theta, \alpha \in \mathbb{R}} M(\theta) + \alpha (\Omega(\theta) - c) = \min_{\theta \in \Theta, \alpha \in \mathbb{R}} \mathcal{L}(\theta, \alpha)$$

- Do **unconstrained** optimization over *both* θ and $\alpha \rightsquigarrow$ derivatives

$$\frac{\partial \mathcal{L}}{\partial \alpha} = \Omega(\theta) - c$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial M}{\partial \theta} + \alpha \frac{\partial \Omega}{\partial \theta}$$

- Set derivatives to zero and solve for θ^*, α^*
- One extra unknown but also one extra equation!

Lagrange multipliers: Constraints \longleftrightarrow Penalties

- Set to zero:

$$\begin{aligned}\Omega(\theta^*) &= c \\ \frac{\partial}{\partial \theta} M(\theta^*) &= -\alpha^* \frac{\partial}{\partial \theta} \Omega(\theta^*)\end{aligned}$$

- First equation is the constraint again
 - So we always get: $\mathcal{L}(\theta^*, \alpha^*) = M(\theta^*)$

Lagrange multipliers: Constraints \longleftrightarrow Penalties

- Set to zero:

$$\begin{aligned}\Omega(\theta^*) &= c \\ \frac{\partial}{\partial \theta} M(\theta^*) &= -\alpha^* \frac{\partial}{\partial \theta} \Omega(\theta^*)\end{aligned}$$

- First equation is the constraint again
 - So we always get: $\mathcal{L}(\theta^*, \alpha^*) = M(\theta^*)$
- Second equation involves *both* θ^* and α^*
- θ^* will **not** be the same as the **unconstrained** optimum θ
 - ...unless the *unconstrained* optimum satisfies the constraint already, in which case $\alpha^* = 0$
- Solving this system might be easy or hard, but the solution does give the **constrained** optimum

Lagrange multipliers are prices

- Changing the constraint level c changes θ^* and $M(\theta)$:

$$\begin{aligned}\frac{\partial M(\theta^*)}{\partial c} &= \frac{\partial \mathcal{L}(\theta^*, \alpha^*)}{\partial c} \\ &= \frac{\partial}{\partial c} \left(M(\theta^*) + \alpha^* (\Omega(\theta^*) - c) \right) \\ &= \left[\frac{\partial}{\partial \theta} M(\theta^*) + \alpha^* \frac{\partial}{\partial \theta} \Omega(\theta^*) \right] \frac{\partial \theta^*}{\partial c} + [\Omega(\theta^*) - c] \frac{\partial \alpha^*}{\partial c} - \alpha^* \\ &= [0] \frac{\partial \theta^*}{\partial c} + [0] \frac{\partial \alpha^*}{\partial c} - \alpha^* \\ &= -\alpha^*\end{aligned}$$

- α^* = Rate at which the optimal *value* improves as the constraint is relaxed
- α^* = How much would you pay for a (marginal) change in the level of the constraint, your **shadow price** for that constraint

Lagrange multipliers vs Penalties

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^* \Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem
- The penalty factor α corresponds to the constraint level c
- "*A fine is a price*"

Moral

Lagrange multipliers turns
constrained optimization
into *penalized* optimization

Lagrange multipliers vs Penalties vs Priors

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*\Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem...

Lagrange multipliers vs Penalties vs Priors

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*\Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem...but don't forget: $M(\theta)$ is typically an *empirical risk* that, we saw, can be thought as a *negative log-likelihood* for some probabilistic model!

Lagrange multipliers vs Penalties vs Priors

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*\Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem...but don't forget: $M(\theta)$ is typically an *empirical risk* that, we saw, can be thought as a *negative log-likelihood* for some probabilistic model!
- Consequence:**
 - At least *formally*, $M(\theta) + \alpha \Omega(\theta)$ can be interpreted from a **Bayesian modeling** perspective as the numerator of a *negative log-posterior* with *negative log-prior* over θ equal to $\alpha \Omega(\theta)$

Lagrange multipliers vs Penalties vs Priors

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*\Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem...but don't forget: $M(\theta)$ is typically an *empirical risk* that, we saw, can be thought as a *negative log-likelihood* for some probabilistic model!
- Consequence:**
 - At least *formally*, $M(\theta) + \alpha \Omega(\theta)$ can be interpreted from a **Bayesian modeling** perspective as the numerator of a *negative log-posterior* with *negative log-prior* over θ equal to $\alpha \Omega(\theta)$
 - Here α is an *hyperparameter* that we may tune (e.g. [Empirical Bayes](#)) or treat [**hierarchically**](#) by selecting a suitable **hyperprior**.

Lagrange multipliers vs Penalties vs Priors

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*\Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem...but don't forget: $M(\theta)$ is typically an *empirical risk* that, we saw, can be thought as a *negative log-likelihood* for some probabilistic model!
- Consequence:**
 - At least *formally*, $M(\theta) + \alpha \Omega(\theta)$ can be interpreted from a **Bayesian modeling** perspective as the numerator of a *negative log-posterior* with *negative log-prior* over θ equal to $\alpha \Omega(\theta)$
 - Here α is an *hyperparameter* that we may tune (e.g. **Empirical Bayes**) or treat **hierarchically** by selecting a suitable **hyperprior**.
 - This interpretation also suggests that, instead of just focusing on *optimization* to get the **MAP** (*maximum a posteriori*) for θ , we may move to **simulation** (e.g. **MCMC methods**) to explore the posterior surface and get a different/better value for θ (+ *uncertainty quantification* too!)

Lagrange multipliers vs Penalties vs Priors

- Once we know α^* , we just do the optimization with an extra term:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*(\Omega(\theta) - c) \\ &= \operatorname{argmin}_{\theta \in \Theta} M(\theta) + \alpha^*\Omega(\theta)\end{aligned}$$

- This is a **penalized** optimization problem...but don't forget: $M(\theta)$ is typically an *empirical risk* that, we saw, can be thought as a *negative log-likelihood* for some probabilistic model!

- Consequence:**

- At least *formally*, $M(\theta) + \alpha \Omega(\theta)$ can be interpreted from a **Bayesian modeling** perspective as the numerator of a *negative log-posterior* with *negative log-prior* over θ equal to $\alpha \Omega(\theta)$
- Here α is an *hyperparameter* that we may tune (e.g. [Empirical Bayes](#)) or treat [hierarchically](#) by selecting a suitable **hyperprior**.
- This interpretation also suggests that, instead of just focusing on *optimization* to get the **MAP** (*maximum a posteriori*) for θ , we may move to **simulation** (e.g. [MCMC methods](#)) to explore the posterior surface and get a different/better value for θ (+ *uncertainty quantification* too!)
- Clearly here we are using the Bayesian prior-to-posterior machinery *pragmatically* as a algorithm-factory: the prior is a tool, not a quantification of our beliefs!

Many constraints

- For multiple equality constraints

$$\Omega_1(\theta) = c_1 \quad \Omega_2(\theta) = c_2 \quad \dots \quad \Omega_n(\theta) = c_n$$

we add n Lagrange multipliers:

$$\mathcal{L}(\theta, \alpha_1, \dots, \alpha_n) = M(\theta) + \sum_{j=1}^n \alpha_j (\Omega_j(\theta) - c_j)$$

- Each constraint equation gets recovered when we take the derivative w.r.t. that multiplier
- Each multiplier tells us our shadow price for loosening each constraint
- Equivalently: adding many *penalty* terms

Inequality constraints

- What if the constraint is that $\Omega(\theta) \leq c$?
- We add on the Lagrange multiplier *anyway*, as though it were an equality

Inequality constraints

- What if the constraint is that $\Omega(\theta) \leq c$?
- We add on the Lagrange multiplier *anyway*, as though it were an equality
- **Case 1:** the global, **unconstrained** optimum obeys the constraint
 - The constraint does not **bind** or **bite**
 - We should get $\alpha^* = 0$

Inequality constraints

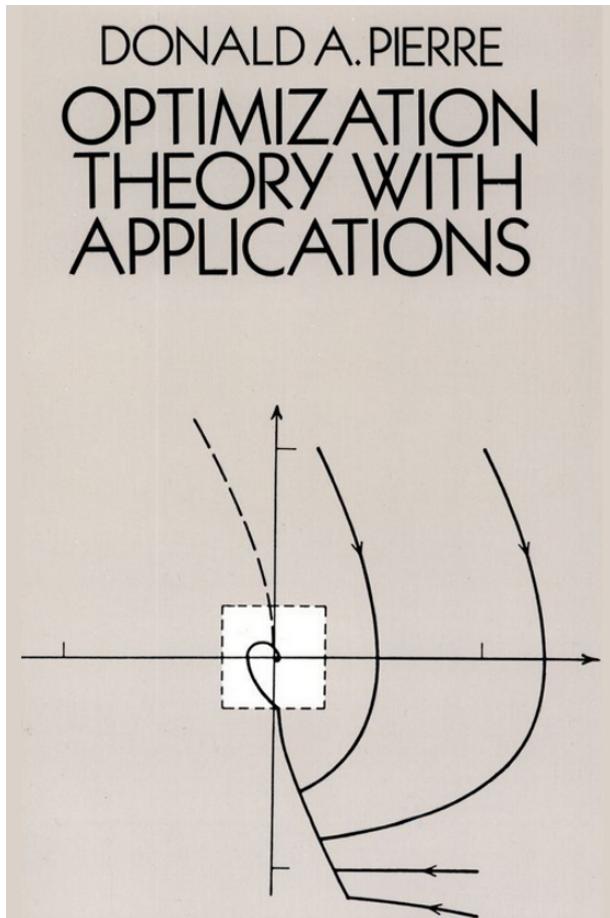
- What if the constraint is that $\Omega(\theta) \leq c$?
- We add on the Lagrange multiplier *anyway*, as though it were an equality
- **Case 1:** the global, **unconstrained** optimum obeys the constraint
 - The constraint does not **bind** or **bite**
 - We should get $\alpha^* = 0$
- **Case 2:** the **unconstrained** optimum is *outside* in constrained feasible set
 - The constraint **binds**, or is **binding**, or **bites**
 - The constrained optimum θ^* is a point where $\Omega(\theta) = c$
 - We'll get $\alpha^* \neq 0$
 - Ignoring some subtleties, this is **Karush-Kuhn-Tucker theorem**

Inequality constraints

- What if the constraint is that $\Omega(\theta) \leq c$?
- We add on the Lagrange multiplier *anyway*, as though it were an equality
- **Case 1:** the global, **unconstrained** optimum obeys the constraint
 - The constraint does not **bind** or **bite**
 - We should get $\alpha^* = 0$
- **Case 2:** the **unconstrained** optimum is *outside* in constrained feasible set
 - The constraint **binds**, or is **binding**, or **bites**
 - The constrained optimum θ^* is a point where $\Omega(\theta) = c$
 - We'll get $\alpha^* \neq 0$
 - Ignoring some subtleties, this is **Karush-Kuhn-Tucker theorem**

Inequality constraints: Remark on KKT conditions

...they were **not** always *this* big...



1st edition
2nd edition
Copyright © 1969, 1986 by Donald A. Pierre.
All rights reserved under Pan American and International
Copyright Conventions.

PREFACE TO THE DOVER EDITION

This Dover edition, first published in 1986, is an unabridged and corrected republication of the work first published by John Wiley & Sons, Inc., New York, 1969. The author has written a new Preface for this edition.

Manufactured in the United States of America
Dover Publications, Inc., 31 East 2nd Street, Mineola,
N.Y. 11501

Library of Congress Cataloging-in-Publication Data

Pierre, Donald A.
Optimization theory with applications.
Reprint. Originally published: New York : Wiley, c1969.
With new pref.
Includes bibliographies and index.
1. Mathematical optimization. 2. Programming
(Mathematics)
I. Title.
QA402.5.P5 1986 519 86-6278
ISBN 0-486-65205-X

The advances in computing capabilities during the past two decades have been impressive. Many types of optimization algorithms are routinely available to the user of a computer center, and small software packages for personal computers are widespread. The understanding of the theory behind these computer routines is important for the user who wishes to gain the most from them or who wishes to embed them as part of a larger system. The strength of this book lies in the breadth of optimization concepts that are presented.

In looking over the contents of the book, my main regret is that I had not included enough on the first- and second-order Kuhn-Tucker conditions which have come to play a major role in many nonlinear programming algorithms. To compensate for this, the major features of the Kuhn-Tucker conditions are presented near the end of this preface; the presentation is

Summing up on constraints and Lagrange multipliers

- Equality constraints act like penalties
 - Penalty factor $\alpha \rightsquigarrow$ Lagrange multiplier enforcing the constraint

Summing up on constraints and Lagrange multipliers

- Equality constraints act like penalties
 - Penalty factor $\alpha \rightsquigarrow$ Lagrange multiplier enforcing the constraint
- Loosening the constraints \rightsquigarrow weakening the penalties

Summing up on constraints and Lagrange multipliers

- Equality constraints act like penalties
 - Penalty factor $\alpha \rightsquigarrow$ Lagrange multiplier enforcing the constraint
- Loosening the constraints \rightsquigarrow weakening the penalties
- Inequality constraints, like $\Omega(\theta) \leq c$, get treated the same way
 - Some multipliers/penalty factors might be 0 if those constraints do **not** bite

Summing up on constraints and Lagrange multipliers

- Equality constraints act like penalties
 - Penalty factor $\alpha \rightsquigarrow$ Lagrange multiplier enforcing the constraint
- Loosening the constraints \rightsquigarrow weakening the penalties
- Inequality constraints, like $\Omega(\theta) \leq c$, get treated the same way
 - Some multipliers/penalty factors might be 0 if those constraints do **not** bite
- Penalties and constraints are different ways of looking at the same thing

Mathematical programming

- Optimization under constraints is called **mathematical programming**
 - The name goes back to the 1930s and is older than "computer programming"
- **Linear programming:** optimize a linear objective function under linear constraints
 - Basically invented by Kantorovich for economic planning in the USSR in the 1930s, re-invented in the West for logistics and decision support in WWII (*Operations Research*), then adapted to corporate decision making, financial portfolio allocation, etc.
- **Convex programming:** optimize a convex function under convex constraints
 - Meaning: take any two points in the feasible set; every point in between them is also in the feasible set
 - Includes linear programming as a special case
- There are efficient (polynomial-time) algorithms for convex programming problems
 - Finding an ϵ -approximate optimum over p variables with k constraints takes time $O((k + p)^{3/2} p^2 \log 1/\epsilon)$

Back to biz!

Learning classifiers by maximizing the margin

- Only consider "linear machines" $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$ with suitably scaled coefficients such that:
 - **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
 - **Subject to:** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$
-

Lagrangian

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 - \sum_i \alpha_i \left\{ y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) - 1 \right\}$$

Learning classifiers by maximizing the margin

- Only consider "linear machines" $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$ with suitably scaled coefficients such that:
 - **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
 - **Subject to:** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$
-

Lagrangian

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 - \sum_i \alpha_i \left\{ y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) - 1 \right\}$$

- Get the *KKT conditions* by setting the gradient of \mathcal{L} w.r.t. the *primal* variables $(\beta_0, \boldsymbol{\beta})$ to zero and then writing the *complementary conditions*

$$\begin{aligned} \nabla_{\boldsymbol{\beta}} \mathcal{L} = \boldsymbol{\beta} - \sum_i \alpha_i y_i \mathbf{x}_i &= 0 &\rightsquigarrow \boldsymbol{\beta}_{\text{svm}} &= \sum_i \alpha_i y_i \mathbf{x}_i \\ \nabla_{\beta_0} \mathcal{L} = - \sum_i \alpha_i y_i &= 0 &\rightsquigarrow \sum_i \alpha_i y_i &= 0 \end{aligned}$$

Learning classifiers by maximizing the margin

- Only consider "linear machines" $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$ with suitably scaled coefficients such that:
 - **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
 - **Subject to:** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$
-

Lagrangian

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 - \sum_i \alpha_i \left\{ y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) - 1 \right\}$$

- Get the *KKT conditions* by setting the gradient of \mathcal{L} w.r.t. the *primal* variables $(\beta_0, \boldsymbol{\beta})$ to zero and then writing the *complementary conditions*

$$\begin{aligned} \nabla_{\boldsymbol{\beta}} \mathcal{L} = \boldsymbol{\beta} - \sum_i \alpha_i y_i \mathbf{x}_i &= 0 &\rightsquigarrow \boldsymbol{\beta}_{\text{svm}} &= \sum_i \alpha_i y_i \mathbf{x}_i \\ \nabla_{\beta_0} \mathcal{L} = -\sum_i \alpha_i y_i &= 0 &\rightsquigarrow \sum_i \alpha_i y_i &= 0 \end{aligned}$$

- If we insert $\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i$ into \mathcal{L} , and use the fact that $\sum_i \alpha_i y_i = 0$, we get the **dual...**

Learning classifiers by maximizing the margin

Lagrangian

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 - \sum_i \alpha_i \left\{ y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) - 1 \right\}$$

KKT conditions

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

Dual

$$F(\boldsymbol{\alpha}) = \underbrace{\frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2 - \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle}_{-\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle} - \underbrace{\sum_i \alpha_i y_i \beta_0 + \sum_i \alpha_i}_{= 0} = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Learning classifiers by maximizing the margin

Lagrangian

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 - \sum_i \alpha_i \left\{ y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) - 1 \right\}$$

KKT conditions

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

Dual

$$F(\boldsymbol{\alpha}) = \underbrace{\frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2 - \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle}_{-\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle} - \underbrace{\sum_i \alpha_i y_i \beta_0 + \sum_i \alpha_i}_{= 0} = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\implies \max_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}) \quad \text{subject to} \quad \alpha_i \geq 0 \quad \forall i \in \{1, \dots, n\} \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

Learning classifiers by maximizing the margin

Lagrangian

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 - \sum_i \alpha_i \left\{ y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) - 1 \right\}$$

KKT conditions

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

Dual

$$\max_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}) \quad \text{subject to} \quad \alpha_i \geq 0 \quad \forall i \in \{1, \dots, n\} \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

- This is a **quadratic program** \rightsquigarrow can be solved quickly!
- We do **not** need the \mathbf{x}_i 's: $F(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
We **only** need to store/know the **inner product** (a.k.a. *similarity measure!*) between pairs of the *original* feature vectors $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = k(\mathbf{x}_i, \mathbf{x}_j)$ \rightsquigarrow *kernelization!*

Consequences of margin maximization: algorithmic

- The new objective function, $F(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$ is a convex (quadratic) function of the α s
- All the constraints on α are convex (linear)
- Convex objective function + convex constraints \Rightarrow unique global minimum with no local minima
- Convex objective + convex constraints = **convex programming** \rightsquigarrow good algorithms which get within ϵ of the optimum in $O(n^{3.5} \log 1/\epsilon)$ time
- Quadratic objective + linear constraints = **quadratic programming**, which can be solved even more efficiently (though not dramatically so)

Moral

Maximizing the margin is algorithmically tractable

Consequences of margin maximization: "support vectors"

- One of our constraints is that for each i

$$y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$$

- This is an inequality constraint
- **Remember:** we can enforce inequality constraints with Lagrange multipliers, but they're only $\neq 0$ if the constraint "binds", so \geq is really =

Consequences of margin maximization: "support vectors"

- One of our constraints is that for each i

$$y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$$

- This is an inequality constraint
- **Remember:** we can enforce inequality constraints with Lagrange multipliers, but they're only $\neq 0$ if the constraint "binds", so \geq is really =
- The multiplier enforcing this constraint for data point i is α_i

Consequences of margin maximization: "support vectors"

- One of our constraints is that for each i

$$y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$$

- This is an inequality constraint
- **Remember:** we can enforce inequality constraints with Lagrange multipliers, but they're only $\neq 0$ if the constraint "binds", so \geq is really =
- The multiplier enforcing this constraint for data point i is α_i
- **Consequence:** $\alpha_i \neq 0$ only for the training points where the margin = 1 \rightsquigarrow active constraint
- If $\alpha_i \neq 0$ then \mathbf{x}_i is a **support vector**
- Generally, most data points won't be support vectors and will have $\alpha_i = 0$, i.e., the solution will be **sparse** (w.r.t. the training data)

Consequences of margin maximization: "support vectors"

- One of our constraints is that for each i

$$y_i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$$

- This is an inequality constraint
- **Remember:** we can enforce inequality constraints with Lagrange multipliers, but they're only $\neq 0$ if the constraint "binds", so \geq is really =
- The multiplier enforcing this constraint for data point i is α_i
- **Consequence:** $\alpha_i \neq 0$ only for the training points where the margin = 1 \rightsquigarrow active constraint
- If $\alpha_i \neq 0$ then \mathbf{x}_i is a **support vector**
- Generally, most data points won't be support vectors and will have $\alpha_i = 0$, i.e., the solution will be **sparse** (w.r.t. the training data)

Consequences of margin maximization: the solution

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

- The optimal hyperplane can be written as

$$f_{\text{svm}}(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}_{\text{svm}} = \beta_0 + \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle = \beta_0 + \sum_{\substack{s : \mathbf{x}_s \\ \text{support point}}} \alpha_s y_s \langle \mathbf{x}_s, \mathbf{x} \rangle$$

Consequences of margin maximization: the solution

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

- The optimal hyperplane can be written as

$$f_{\text{svm}}(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}_{\text{svm}} = \beta_0 + \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle = \beta_0 + \sum_{\substack{s : \mathbf{x}_s \\ \text{support point}}} \alpha_s y_s \langle \mathbf{x}_s, \mathbf{x} \rangle$$

- Since support vectors lie on the marginal hyperplanes, for *any* support point, say \mathbf{x}_s , we get

$$y_s = \beta_0 + \langle \boldsymbol{\beta}_{\text{svm}}, \mathbf{x}_s \rangle \quad \rightsquigarrow \quad \beta_0 = y_s - \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x}_s \rangle$$

Consequences of margin maximization: the solution

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

- The optimal hyperplane can be written as

$$f_{\text{svm}}(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}_{\text{svm}} = \beta_0 + \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle = \beta_0 + \sum_{\substack{s : \mathbf{x}_s \\ \text{support point}}} \alpha_s y_s \langle \mathbf{x}_s, \mathbf{x} \rangle$$

- Since support vectors lie on the marginal hyperplanes, for *any* support point, say \mathbf{x}_s , we get

$$y_s = \beta_0 + \langle \boldsymbol{\beta}_{\text{svm}}, \mathbf{x}_s \rangle \quad \rightsquigarrow \quad \beta_0 = y_s - \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x}_s \rangle$$

- We can also derive a simple expression for the geometric margin ρ in terms of $\boldsymbol{\alpha}$ by multiplying both side of the last equation by $\alpha_s y_s$ and summing up (over support points)

$$\sum_s \alpha_s y_s \beta_0 = \sum_s \alpha_s y_s^2 - \sum_{i,s} \alpha_i \alpha_s y_i y_s \langle \mathbf{x}_i, \mathbf{x}_s \rangle \quad \Leftrightarrow \quad 0 = \sum_i \alpha_i - \|\boldsymbol{\beta}\|_2^2 \quad \Leftrightarrow \quad \rho^2 \stackrel{\text{def}}{=} \frac{1}{\|\boldsymbol{\beta}\|_2^2} = \frac{1}{\sum_s \alpha_s} \stackrel{\alpha_s > 0}{=} \frac{1}{\|\boldsymbol{\alpha}\|_1}$$

Level 2: Slightly Harder...but Soft...

The Non-Linearly Separable Case

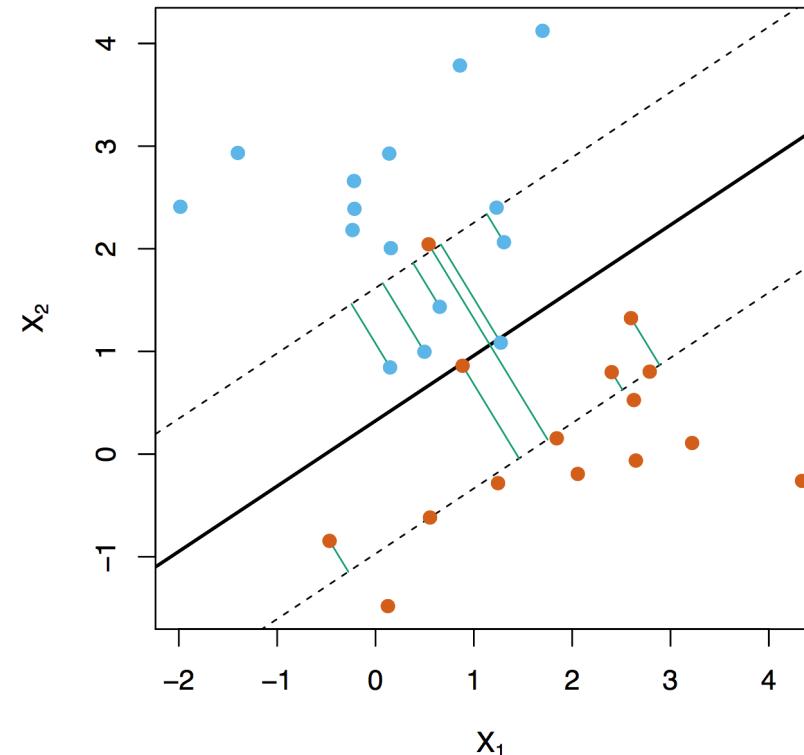
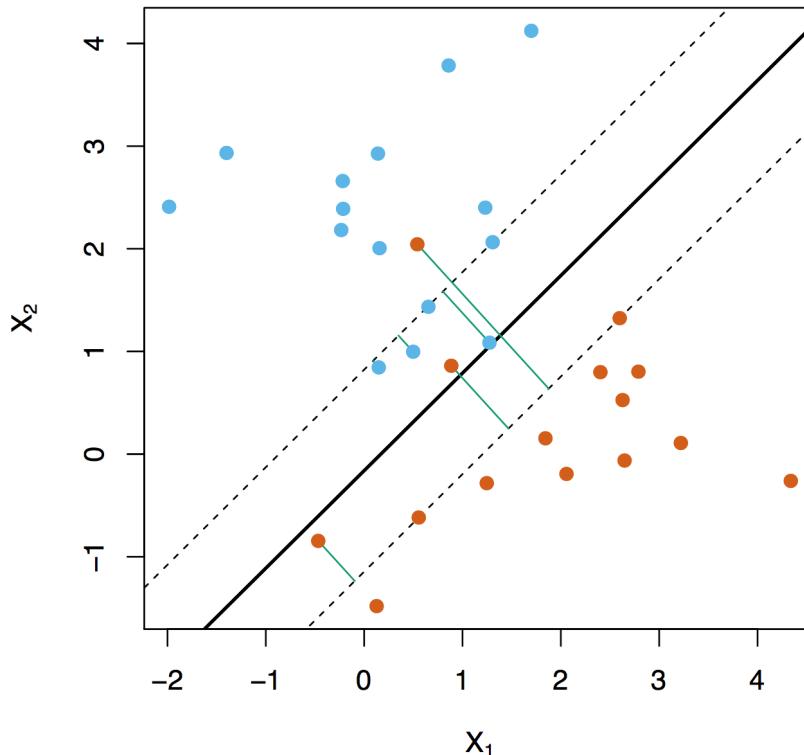
(Soft Margin)

"Soft" margin and slack

- Sometimes we just can't get perfect classification with a plane $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$

"Soft" margin and slack

- Sometimes we just can't get perfect classification with a plane $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$
- **Solution:** allow for margin violations under a *budget constraint* \rightsquigarrow **tuning parameter**



"Soft" margin and slack

- Sometimes we just can't get perfect classification with a plane $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$
 - **Solution:** allow for margin violations under a *budget constraint* \rightsquigarrow **tuning parameter**
-

Before

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
- **Subject to:** hard separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$

"Soft" margin and slack

- Sometimes we just can't get perfect classification with a plane $f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$
 - **Solution:** allow for margin violations under a *budget constraint* \rightsquigarrow **tuning parameter**
-

Before

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
 - **Subject to:** hard separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1$ for each $i \in \{1, \dots, n\}$
-

After

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$
 - **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
-

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
 - **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
-

- The slack $\xi_i > 0$ if \mathbf{x}_i is "beyond" the margin \rightsquigarrow we pay a penalty equal to $\sum_i \xi_i$
- C is the budget for the total amount of violations we allow.

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
- **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
- The slack $\xi_i > 0$ if \mathbf{x}_i is "beyond" the margin \rightsquigarrow we pay a penalty equal to $\sum_i \xi_i$
- C is the budget for the total amount of violations we allow.
- The solution to this new optimization problem has the same shape as before, except now the support set includes any vector on the margin as well as those that violate the margin.

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
- **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$

- The slack $\xi_i > 0$ if \mathbf{x}_i is "beyond" the margin \rightsquigarrow we pay a penalty equal to $\sum_i \xi_i$
- C is the budget for the total amount of violations we allow.
- The solution to this new optimization problem has the same shape as before, except now the support set includes any vector on the margin as well as those that violate the margin.
- The bigger $C \rightsquigarrow$ the larger the support set \rightsquigarrow the more points that have an impact on the solution.

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
- **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$

- The slack $\xi_i > 0$ if \mathbf{x}_i is "beyond" the margin \rightsquigarrow we pay a penalty equal to $\sum_i \xi_i$
- C is the budget for the total amount of violations we allow.
- The solution to this new optimization problem has the same shape as before, except now the support set includes any vector on the margin as well as those that violate the margin.
- The bigger $C \rightsquigarrow$ the larger the support set \rightsquigarrow the more points that have an impact on the solution.
- Hence bigger $C \rightsquigarrow$ more stability \rightsquigarrow lower variance.
- In fact even for separable data, allowing margin violations lets us *regularize* the solution.

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
 - **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
-

- **Notice:** the constraint can be succinctly captured as follows

$$\forall i \in \{1, \dots, n\}, \xi_i \geq 1 - y_i f(\mathbf{x}_i) \rightsquigarrow C = \sum_i \xi_i \geq \sum_i (1 - y_i f(\mathbf{x}_i)) \stackrel{(\heartsuit)}{\Leftrightarrow} \sum_i (1 - y_i f(\mathbf{x}_i))_+ \leq C,$$

where we have introduced the positive part $(u)_+ = \max\{0, u\}$ in (\heartsuit) because we must pay a price $\xi_i > 0 \Leftrightarrow$ there's a violation \Leftrightarrow the margin $y_i f(\mathbf{x}_i)$ is less than 1 $\Leftrightarrow 1 - y_i f(\mathbf{x}_i) > 0$.

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
 - **Subject to:**
 - **soft separability condition** $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
-

- **Notice:** the constraint can be succinctly captured as follows

$$\forall i \in \{1, \dots, n\}, \xi_i \geq 1 - y_i f(\mathbf{x}_i) \rightsquigarrow C = \sum_i \xi_i \geq \sum_i (1 - y_i f(\mathbf{x}_i)) \stackrel{(\heartsuit)}{\Leftrightarrow} \sum_i (1 - y_i f(\mathbf{x}_i))_+ \leq C,$$

where we have introduced the positive part $(u)_+ = \max\{0, u\}$ in (\heartsuit) because we must pay a price $\xi_i > 0$
 \Leftrightarrow there's a violation \Leftrightarrow the margin $y_i f(\mathbf{x}_i)$ is less than 1 $\Leftrightarrow 1 - y_i f(\mathbf{x}_i) > 0$.

Lagrangian

$$\min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2 + \alpha \sum_i (1 - y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle))_+ \rightsquigarrow \text{multiply by } \lambda = \frac{1}{\alpha} \text{ (and forget the const. } 1/2)$$

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
 - **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
-

Lagrangian

$$\min_{(\beta_0, \beta)} \sum_i (1 - y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle))_+ + \lambda \|\beta\|_2^2 \stackrel{\mathcal{F} = \{\text{Linear in } \mathbf{x}\}}{\rightsquigarrow} \min_{f \in \mathcal{F}} \widehat{R}_{\mathbb{H}}(f) + \lambda \text{pen}(f)$$

More bias/Less variance \rightsquigarrow Large C \rightsquigarrow Small α \rightsquigarrow Large λ \rightsquigarrow Heavily regularised

"Soft" margin and slack

- **Maximize** the margin $M(f) \rightsquigarrow \min_{(\beta_0, \beta)} \frac{1}{2} \|\beta\|_2^2$
 - **Subject to:**
 - **soft** separability condition $\rightsquigarrow y_i(\beta_0 + \langle \beta, \mathbf{x}_i \rangle) \geq 1 - \xi_i$ for each $i \in \{1, \dots, n\}$
 - **budget constraint** $\rightsquigarrow \xi_i \geq 0$ and $\sum_i \xi_i = C$
-

Lagrangian

$$\min_{(\beta_0, \beta)} \sum_i (1 - y_i (\beta_0 + \langle \beta, \mathbf{x}_i \rangle))_+ + \lambda \|\beta\|_2^2 \stackrel{\mathcal{F} = \{\text{Linear in } \mathbf{x}\}}{\rightsquigarrow} \min_{f \in \mathcal{F}} \widehat{R}_{\mathbb{H}}(f) + \lambda \text{pen}(f)$$

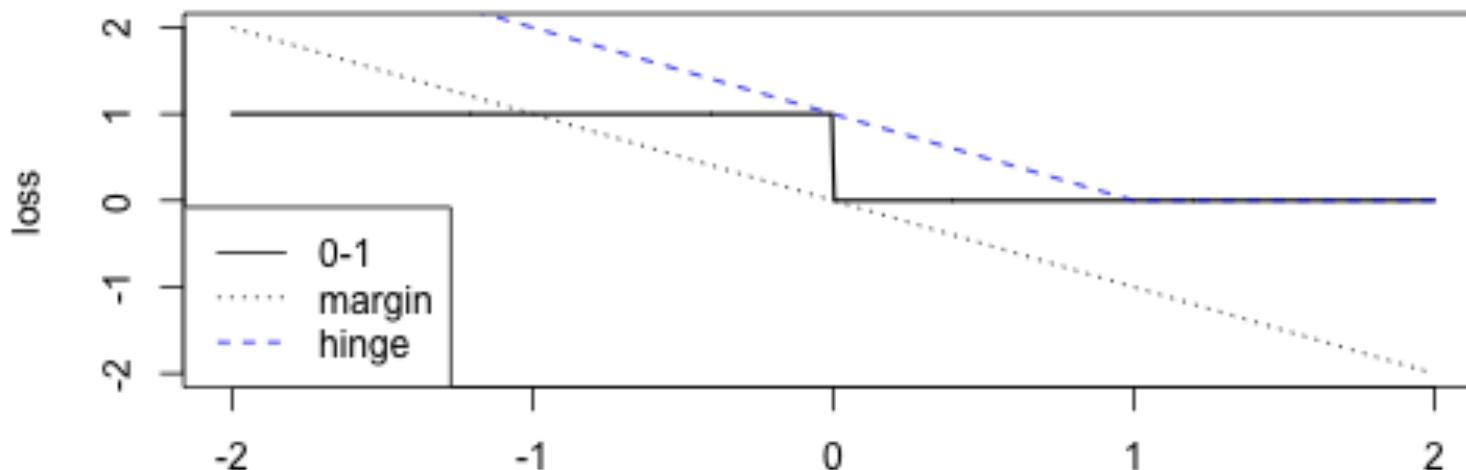
More bias/Less variance \rightsquigarrow Large C \rightsquigarrow Small α \rightsquigarrow Large λ \rightsquigarrow Heavily regularised

Conclusion

As promised, from purely geometric principles
we obtained (almost) for free two fundamental analytic learning tools:
Surrogate Loss + Regularization

Hinge loss/risk

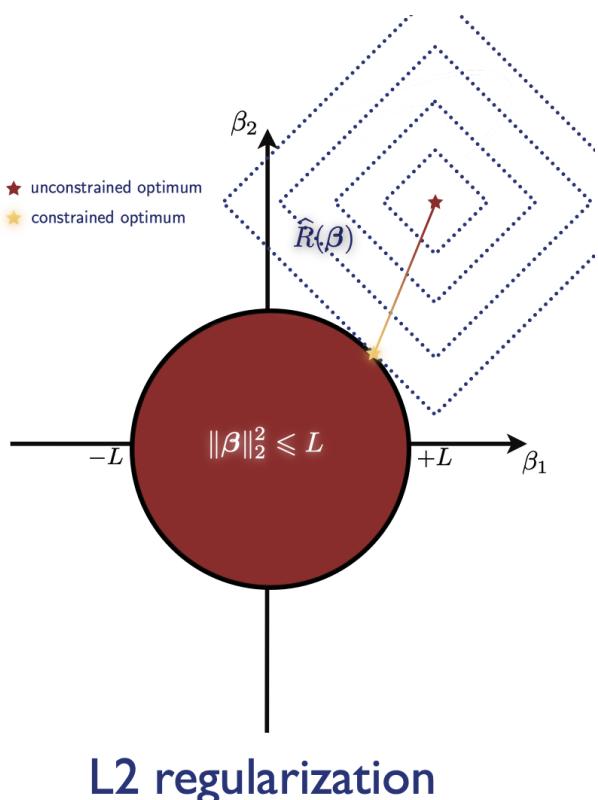
$$L(y, f(x)) = (1 - y f(x))_+ \equiv \max \{0, 1 - y f(x)\}$$



- **Remember:** the 0-1 loss is always \leq hinge loss, so if we can bound the hinge loss we will also bound the miss-classification/accuracy risk!

SVM: Remarks

Penalty \leftrightarrow Constraint



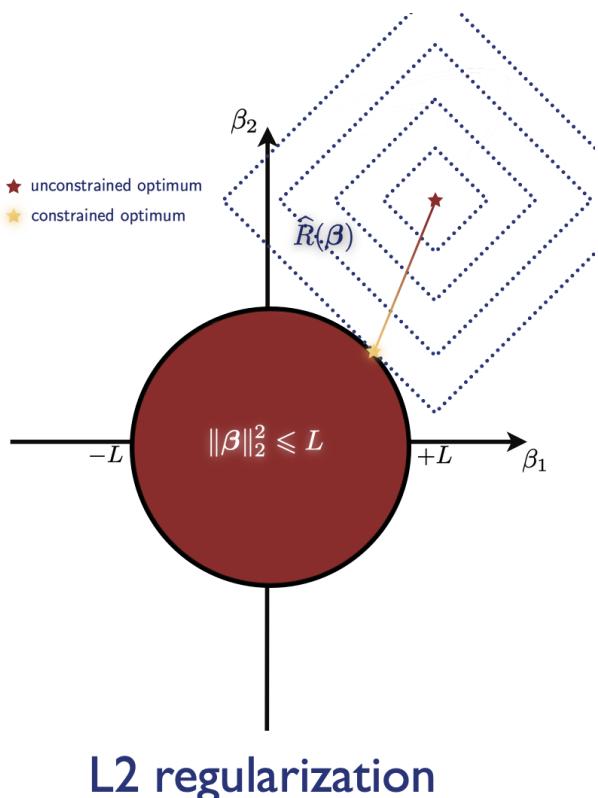
- Hinge and penalty are both **convex** in $\beta \rightsquigarrow$ for each λ we can find an L s.t. the following constrained problem has the same solution

$$\min_{(\beta_0, \beta)} \hat{R}_H(\beta_0, \beta) \quad \text{subject to} \quad \|\beta\|_2^2 \leq L$$

- Typically the bias/intercept β_0 is *not* penalized.
- Typically we do **not** solve the optimization problem in this form, but we go for the *primal-dual* way detailed before.
- This rewriting hints to a plethora of variations:
 - **Change penalty:** $\|\beta\|_2 = L^2 \rightsquigarrow L^1 = \|\beta\|_1$ to force *sparsity* in β (e.g. **LASSO**, more soon)
This is known as L^1 -SVM (**penalizedSVM**)
 - *Pro:* it can be written as a **linear** programming problem
 - *Con:* **no** longer a maximum margin classifier!

SVM: Remarks

Penalty \leftrightarrow Constraint



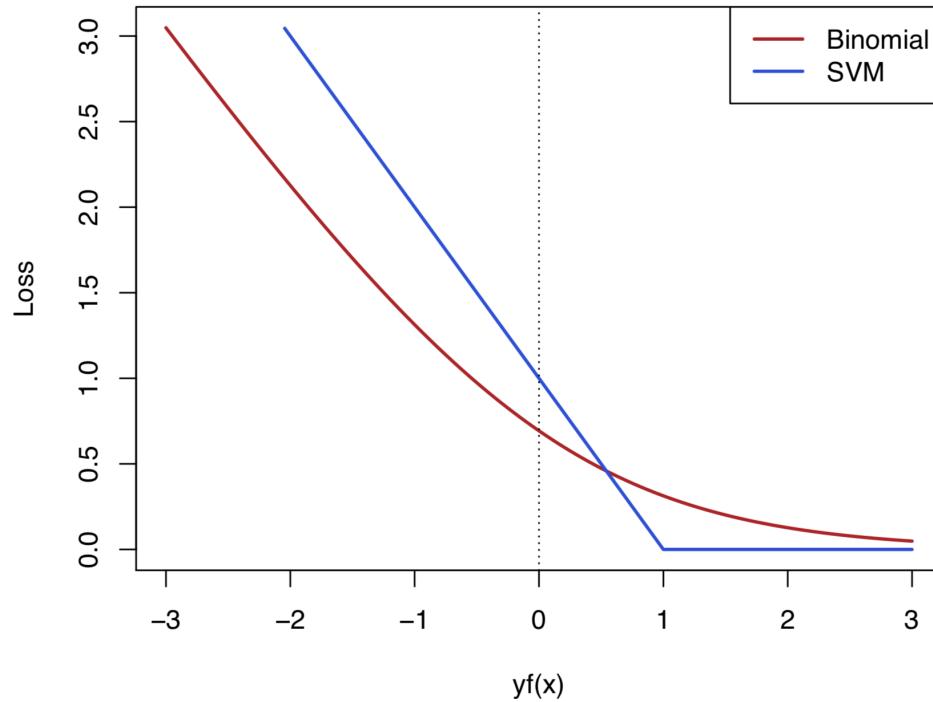
- Hinge and penalty are both **convex** in $\beta \rightsquigarrow$ for each λ we can find an L s.t. the following constrained problem has the same solution

$$\min_{(\beta_0, \beta)} \hat{R}_H(\beta_0, \beta) \quad \text{subject to} \quad \|\beta\|_2^2 \leq L$$

- Typically the bias/intercept β_0 is *not* penalized.
- Typically we do **not** solve the optimization problem in this form, but we go for the *primal-dual* way detailed before.
- This rewriting hints to a plethora of variations:
 - **Change penalty:** $\|\beta\|_2 = L^2 \rightsquigarrow L^1 = \|\beta\|_1$ to force *sparsity* in β (e.g. **LASSO**, more soon)
 - **Change loss:** to handle different tasks (e.g. Regression) \rightsquigarrow if done properly the *primal-dual* route may still be feasible
 - **Change optimizer:** to handle big data, go **SGD** by (numerically) differentiate it

SVM: Remarks

Penalty \longleftrightarrow Constraint



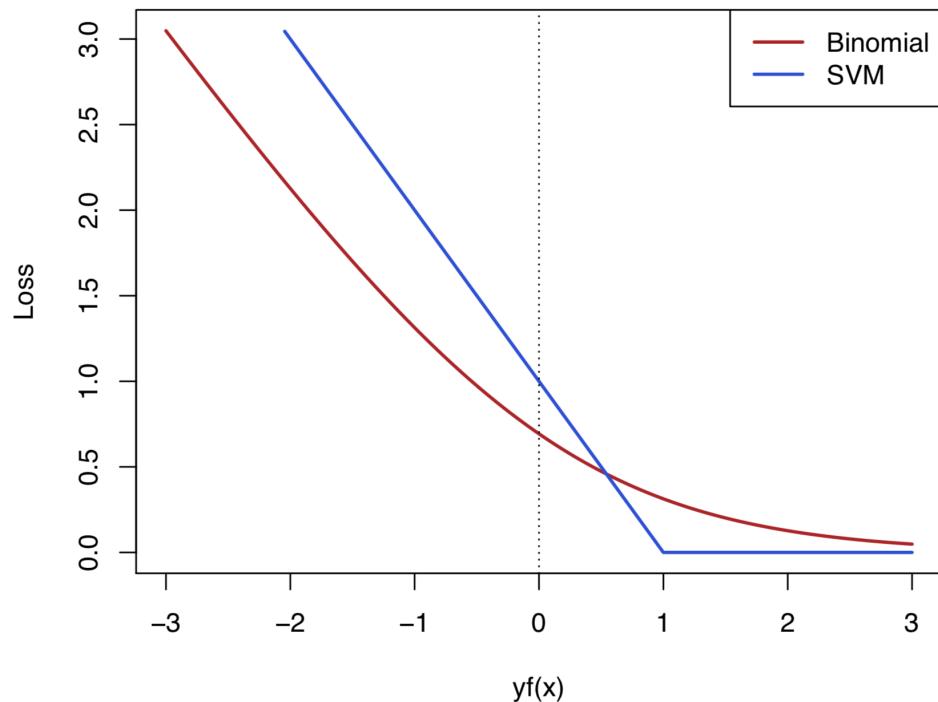
- *Penalized Empirical Risk* methods are ubiquitous: given two solutions $\hat{\beta}_1$ and $\hat{\beta}_2$ with the same (empirical) risk, we prefer the more parsimonious of the two (as measured by the selected penalty/norm)
- For example, we know the **logistic-loss** needs some help (even under linear separability)
- We hint at *early-stopping* as a *regularizer*: the number of iterations is the tuning parameter to select via cross-validation, for example
- Another option is clearly to penalize the **logistic deviance/loss**, here for $Y \in \{-1, +1\}$

$$\min_{(\beta_0, \beta)} \sum_{i=1}^n \log [1 + e^{-y_i (\beta_0 + \langle \mathbf{x}_i, \beta \rangle)}] + \lambda \|\beta\|_2^2.$$

Under some conditions, *early stopping* $\approx L^2$
(see **ESL**, Section 10.12 or 16.2.3)

SVM: Remarks

Penalty \longleftrightarrow Constraint



$$\min_{(\beta_0, \beta)} \sum_{i=1}^n \log [1 + e^{-y_i (\beta_0 + \langle \mathbf{x}_i, \beta \rangle)}] + \lambda \|\beta\|_2^2.$$

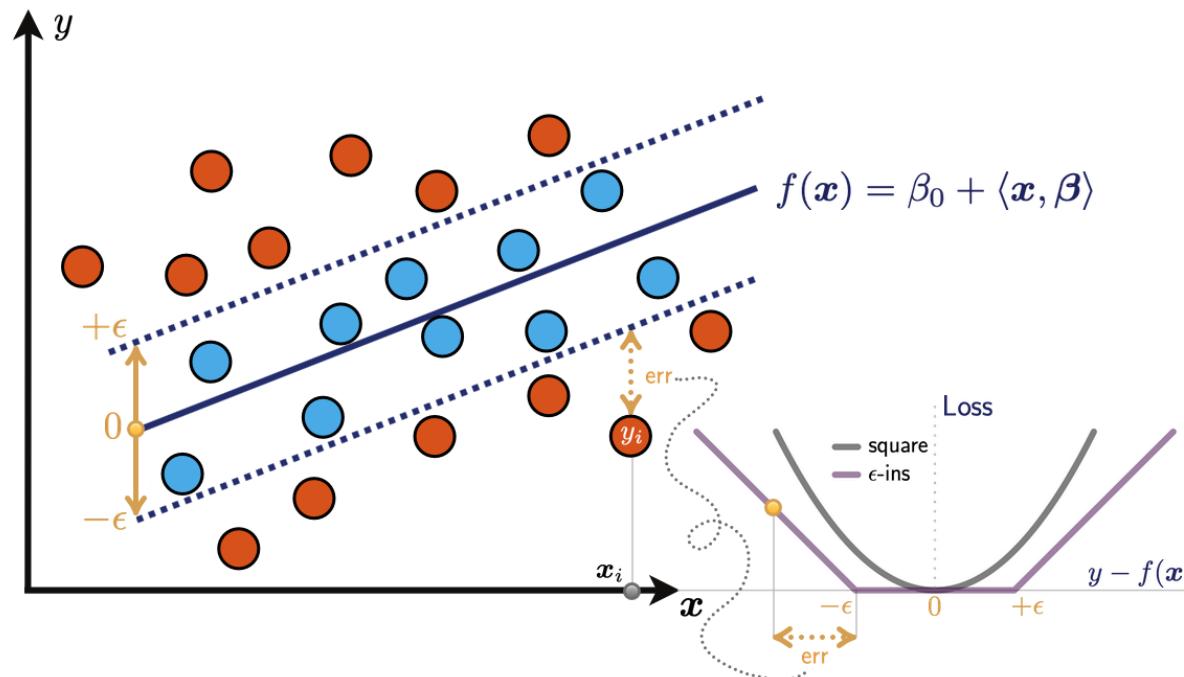
- *Hinge*: elbow at $+1$ (vs) *Logistic*: smooth bend
 \leadsto the logistic-solution involves **all** the data:
the weights $f_i(1 - f_i)$ fade smoothly with
distance from the decision boundary as
opposed to the binary nature of support points.
- **Rosset et al. (2004)**: for $\lambda \downarrow 0$, the scaled logistic
coeff's $\beta / \|\beta\|_2$ converges to the SVM solution
under linear separability of the data.
 - Because of the required normalization for
the logistic, the SVM solution is preferable.
 - On the other hand, for non-linearly
separable cases, the logistic solution may
have some advantages since its target is
the *logit* (log-odds) of class probabilities.
- Check the **liblinearR** package.

SVM: Remarks

Support Vector Regression

- **Hope:** by replacing the hinge with a suitable loss, we may experience the beauty of SVM (i.e. easy optimization problem that depends on inner-products *only* + data-sparse solution) also in *regression*...

$$L(y, f(\mathbf{x})) = |y - f(\mathbf{x})|_{\epsilon} = \max \{0, |y - f(\mathbf{x})| - \epsilon\} \quad \epsilon\text{-insensitive loss}$$

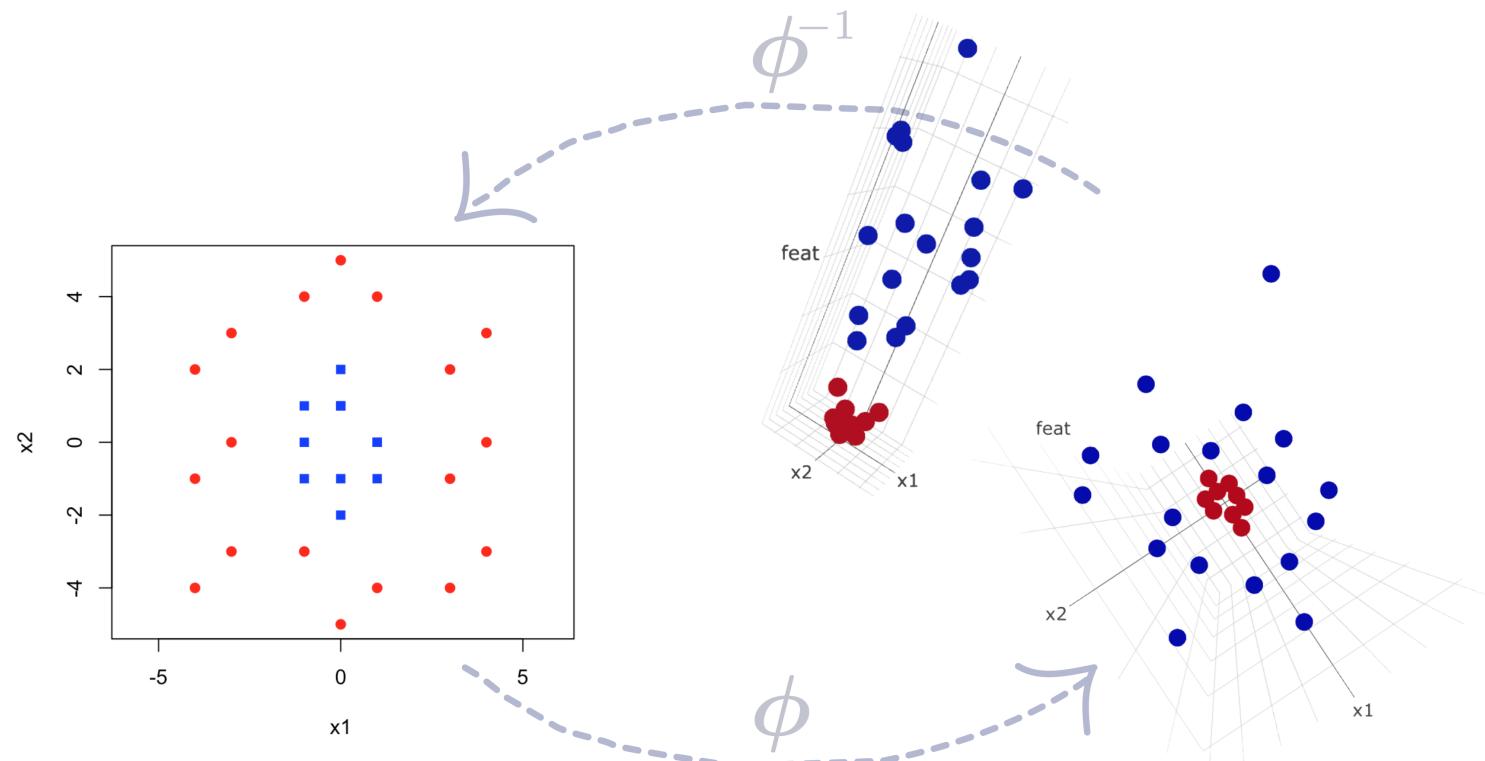


Level 3: Kernelization

From Linear to Nonlinear in a Trick!

Explicit Nonlinearity in Classification

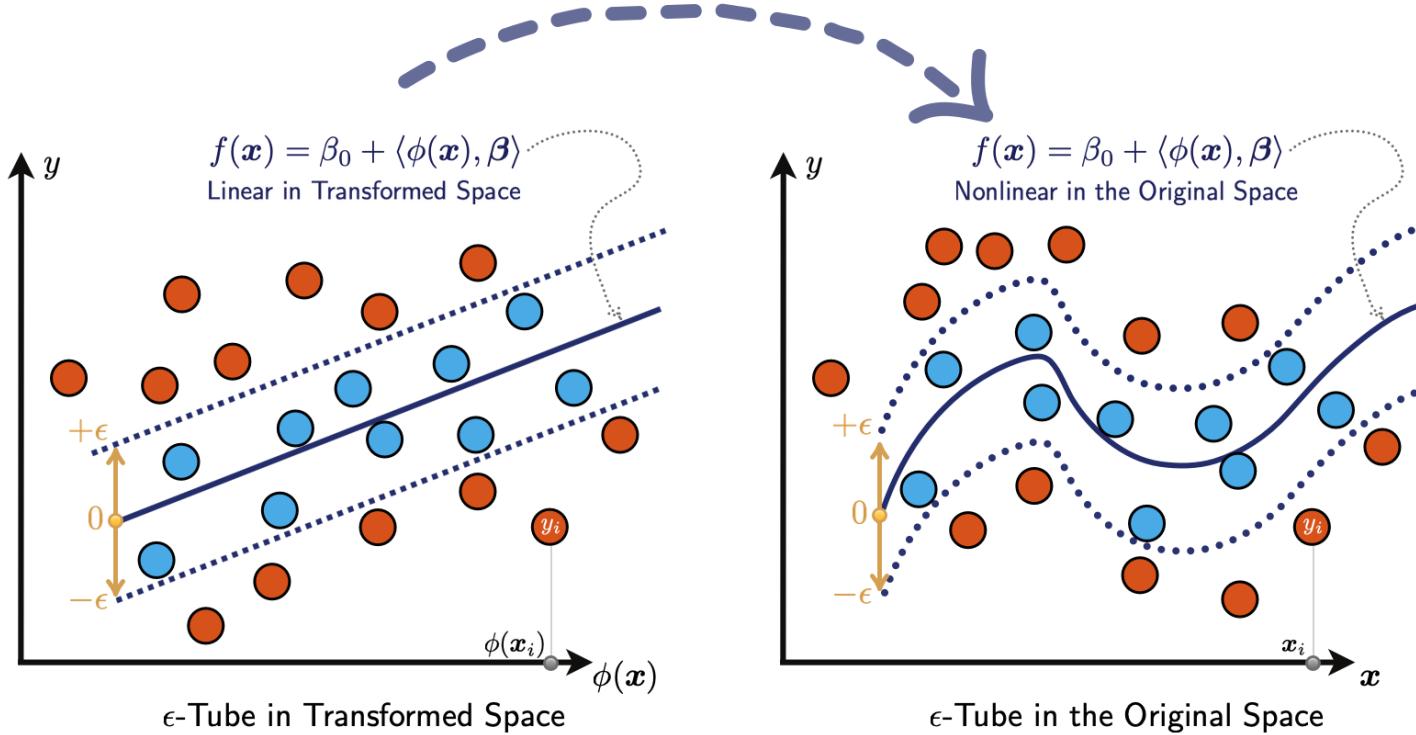
Nonlinearly Separable in the Original Space \iff Linearly Separable in Transformed Space



$$\mathbf{x} = [x_1 \quad x_2]^T \rightsquigarrow \phi(\mathbf{x}) = \boldsymbol{\phi} = [\phi_1 \quad \phi_2 \quad \phi_3]^T = [x_1^2 \quad x_2^2 \quad \sqrt{2}x_1x_2]^T$$

Explicit Nonlinearity in Regression

Support Vector Regression



Back in the SVM-days...

Dual

$$F(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\max_{\boldsymbol{\alpha}} F(\boldsymbol{\alpha}) \quad \text{subject to} \quad \alpha_i \geq 0 \quad \forall i \in \{1, \dots, n\} \quad \text{and} \quad \sum_i \alpha_i y_i = 0$$

Optimal hyperplane

$$\boldsymbol{\beta}_{\text{svm}} = \sum_i \alpha_i y_i \mathbf{x}_i \rightsquigarrow f_{\text{svm}}(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}_{\text{svm}} = \beta_0 + \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle = \beta_0 + \sum_{\substack{s : \mathbf{x}_s \\ \text{support point}}} \alpha_s y_s \langle \mathbf{x}_s, \mathbf{x} \rangle$$

Consequences

Working in the original \mathcal{X} space, we do **not** need the \mathbf{x}_i 's

Training: store **inner products** (*similarity measure*) between pairs of *training data* $\rightsquigarrow \langle \mathbf{x}_i, \mathbf{x}_j \rangle$

Prediction: evaluate **inner products** between the new feature vector \mathbf{x} and the *training data* $\rightsquigarrow \langle \mathbf{x}_i, \mathbf{x} \rangle$

Back in the SVM-days...

Dual (matrix form)

$$F(\boldsymbol{\gamma}) = \mathbf{1}^T \boldsymbol{\alpha} - \frac{1}{2} (\mathbb{Y}\boldsymbol{\alpha})^T \mathbb{K} (\mathbb{Y}\boldsymbol{\alpha})^{\gamma=\mathbb{Y}\boldsymbol{\alpha}} \mathbf{1}^T (\mathbb{Y}^{-1}\boldsymbol{\gamma}) - \frac{1}{2} \boldsymbol{\gamma}^T \mathbb{K} \boldsymbol{\gamma} \quad \text{with} \quad \mathbf{1}^T = [1 \ 1 \ \dots \ 1] \in \mathbb{R}^n \quad \text{and} \quad \mathbb{Y} = \text{diag}(\mathbf{y})$$

$$\max_{\gamma} F(\gamma) \quad \text{subject to} \quad (\mathbb{Y}^{-1}\gamma) \succcurlyeq 0 \quad \text{and} \quad \mathbf{1}^T \gamma = 0$$

Optimal hyperplane (matrix form)

Consequences

Working in the original \mathcal{X} space, we do **not** need the x_i 's

Training: store the $(n \times n)$ **Gramian matrix** $\mathbb{K} = (\mathbb{X} \mathbb{X}^T) \iff \mathbb{K}[i, j] = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$

Prediction: evaluate the **inner product** vector $\mathbb{X} \mathbf{x}$ (on support points only!)

Interlude: $f(\mathbf{x}) = \beta_0 + \langle \mathbf{x}, \boldsymbol{\beta} \rangle = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}$

Support Vector (Classification) + Support Vector (Regression) vs Ordinary Least Squares

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i f(\mathbf{x}_i))_+ \quad \text{vs} \quad \frac{1}{n} \sum_{i=1}^n |y_i - f(\mathbf{x}_i)|_\epsilon \quad \text{vs} \quad \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

- *Computations*: differently from OLS/square-loss, hinge and ϵ -insensitive do **not** allow for a closed form solution, but we are happy with the efficiency of their *primal/dual* quadratic-programming formulation.
 - *Primal/Dual view*: we have seen that the SVM-solution can be more naturally expressed in *dual* form:

This is only apparently in contrast with the familiar closed form solution to the OLS problem. In fact...

$$\boldsymbol{\beta}_{\text{ols}} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y} \quad \rightsquigarrow \quad f_{\text{ols}}(\mathbf{x}) = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}_{\text{ols}} = \beta_0 + \mathbf{y}^T (\mathbb{X} \mathbb{X}^T)^{-1} \mathbb{X} \mathbf{x} = \beta_0 + \boldsymbol{\gamma}_{\text{ols}}^T \mathbb{X} \mathbf{x}$$

where now the optimal (*dual*) coefficients admit a closed form $\gamma_{\text{ols}} = (\mathbb{K}^{-1})^T \mathbf{y} \stackrel{\text{sym}}{=} \mathbb{K}^{-1} \mathbf{y}$ ↪ **kernelizable!**

Toward Nonlinearity...and Beyond!

A) Classic route \longleftrightarrow Explicit transformation

Idea: Embed the problem in a larger one by **explicitly** defining a non-linear transformation $\phi(\cdot)$, the **feature map**, which associates a possibly very high dimensional vector $\phi = \phi(\mathbf{x}) \in \mathbb{R}^q$ to any feature profile $\mathbf{x} \in \mathcal{X}$.

Goal: Engineer the map $\phi(\cdot)$ in such a way that: 1. in the transformed space, linearity is enough to achieve our predictive goals, 2. there is a computationally efficient way to work it out

- *Remark:* generally speaking, for 1. we need some faith; while for 2, we have solid results

Practice: once the *feature map* $\phi(\cdot)$ has been designed, we can simply work **explicitly** in the new covariates space - as in the case of *polynomial/spline regression* - rewriting the optimization problem in an obvious way

- *Remark:* solving it, we still get a *linear* SVM-classifier **but** in the ϕ -transformed space!

$$F(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \phi_i, \phi_j \rangle_{\mathbb{R}^q} \quad \text{and} \quad f_{\text{svm}}(\mathbf{x}) = \beta_0 + \sum_i \alpha_i y_i \langle \phi_i, \phi(\mathbf{x}) \rangle_{\mathbb{R}^q}$$

Analysis: the SVM computations scale linearly in q (although potentially cubic in n) \rightsquigarrow doable but typically **not** efficient for very large q (in the tens of thousands or even millions) \rightsquigarrow can we avoid remapping all the \mathbf{x}_i ?

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Idea: since what we really need is only the inner product $\langle \phi_i, \phi_j \rangle_{\mathbb{R}^q}$, we may try do design $\phi(\cdot)$ so that $\langle \phi_i, \phi_j \rangle_{\mathbb{R}^q}$ admits a simple representation in terms of the original covariate $\{\mathbf{x}_i\}_i \rightsquigarrow$ **no** explicit evaluation of the embeddings $\{\phi_i\}_i$ is needed!

Example 1: let's go back to the initial example, and define the following *feature map* $\phi : \mathbb{R}^2 \mapsto \mathbb{R}^3$:

$$\mathbf{x} = [x_1 \quad x_2]^T \rightsquigarrow \phi(\mathbf{x}) = \boldsymbol{\phi} = [\phi_1 \quad \phi_2 \quad \phi_3]^T = [x_1^2 \quad x_2^2 \quad \sqrt{2}x_1x_2]^T$$

Given two feature profiles \mathbf{x} and \mathbf{y} , let us now evaluate the inner-product:

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathbb{R}^3} &= [x_1^2 \quad x_2^2 \quad \sqrt{2}x_1x_2] \begin{bmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2}y_1y_2 \end{bmatrix} = x_1^2y_1^2 + x_2^2y_2^2 + 2(x_1y_1)(x_2y_2) \\ &= (x_1y_1 + x_2y_2)^2 = \langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{R}^2}^2 = K(\mathbf{x}, \mathbf{y}) \quad (\text{example of polynomial kernel})\end{aligned}$$

Mission accomplished!

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Idea: since what we really need is only the inner product $\langle \phi_i, \phi_j \rangle_{\mathbb{R}^q}$, we may try do design $\phi(\cdot)$ so that $\langle \phi_i, \phi_j \rangle_{\mathbb{R}^q}$ admits a simple representation in terms of the original covariate $\{\mathbf{x}_i\}_i \rightsquigarrow$ **no** explicit evaluation of the embeddings $\{\phi_i\}_i$ is needed!

Example 1: more in general, functions of two variables \mathbf{x} and \mathbf{y} both in $\mathcal{X} \subseteq \mathbb{R}^p$ with the following form

$$K(\mathbf{x}, \mathbf{y}) = (c + \langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{R}^p})^d \quad \text{with} \quad d \in \mathbb{N}, c > 0 \quad (\text{to be tuned via CV})$$

are called **polynomial kernels**, and *implicitly* code for an inner product in a high-dimensional space.

- Specifically, once we expand the power d and express K in terms of *all* monomials $x_1^{k_1} x_2^{k_2} \cdots x_p^{k_p}$, we are left with the problem of determining the number of monomials of degree exactly d in p input variables
- This is equivalent to find the number of ways of adding p non-negative integers to obtain a sum of at most d that can be shown to be equal to $\binom{p+d}{d}$
- In addition, by the **binomial theorem**, we deduce that the weight assigned to the monomial $x_j^{k_j}$ is $\binom{d}{j} c^{d-j}$ for each $j \leq d \rightsquigarrow$ increasing c decreases the weight of k_j , particularly that of k_j 's with larger j .

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 2: Gaussian Kernel/Radial Basis Function Kernel

Important Remark

As long as we can compute inner products, there is **no** restriction on the dimension of the feature space.

Immediate Consequence

The most commonly used kernel corresponds to "projection" into an **infinite** dimensional space...*implicitly*, of course!

Assume for simplicity that $x \in \mathcal{X} = \mathbb{R} \rightsquigarrow d = 1$ and define

$$\phi_j(x) = x^j \times \left(\frac{2^j}{j!}\right)^{\frac{1}{2}} \times e^{-x^2} \quad j \in \{0, 1, 2, \dots\} = \mathbb{N}_0 \rightsquigarrow \boldsymbol{\phi} = [\phi_0(x) \quad \phi_1(x) \quad \phi_2(x) \cdots]^T$$

power part $\downarrow j$: damping factor normalization

Remarks:

- Notice we are defining a *feature map* $\phi : \mathcal{X} \subseteq \mathbb{R} \mapsto \ell^2 \approx \mathbb{R}^\infty$ in a *infinite dimensional sequence space*!
- *Normalization:* $\|\boldsymbol{\phi}\|_{\ell^2}^2 = \sum_{j=0}^{\infty} \phi_j^2(x) = \sum_{j=0}^{\infty} x^{2j} \left(\frac{2^j}{j!}\right) (e^{-x^2})^2 = e^{-2x^2} \sum_{j=0}^{\infty} \frac{(2x)^j}{j!} \stackrel{\text{Taylor}}{=} e^{-2x^2} e^{2x^2} = 1$ (normalized!)

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 2: Gaussian Kernel/Radial Basis Function Kernel

$$\phi_j(x) = x^j \times \left(\frac{2^j}{j!}\right)^{\frac{1}{2}} \times e^{-x^2} \quad j \in \{0, 1, 2, \dots\} = \mathbb{N}_0 \quad \rightsquigarrow \quad \boldsymbol{\phi} = [\phi_0(x) \quad \phi_1(x) \quad \phi_2(x) \cdots]^T$$

- In essence, we are mapping any $x \in \mathbb{R}$ into a norm 1, infinite dimensional vector $\boldsymbol{\phi}$ or, in other words, we are remapping the entire input space $\mathcal{X} = \mathbb{R}$ into the unit ball of the sequence space ℓ^2 .
- As with the *polynomial kernel*, let us now look at the actual shape of the kernel implied by these $\{\phi_j\}_j$. Let x and y be two input values and evaluate their inner product once transformed by our *feature map* ϕ :

$$\langle \phi(x), \phi(y) \rangle_{\ell^2} = e^{-x^2} e^{-y^2} \sum_{j=0}^{\infty} \frac{(2xy)^2}{j!} = e^{-(x^2 - 2xy + y^2)} = e^{-(x-y)^2} = K(x, y) \quad (\text{Gaussian kernel})$$

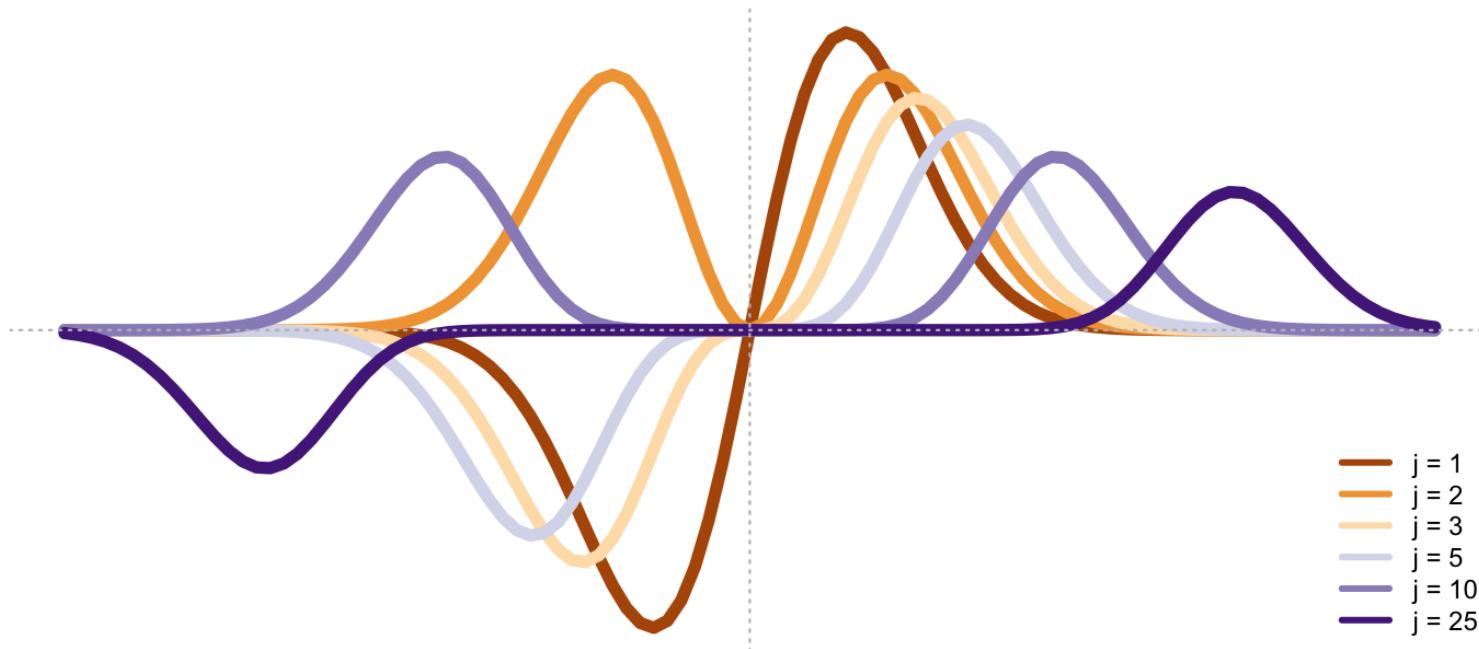
- If $\mathbf{x} \in \mathbb{R}^p$, we can define a (parametrized) Gaussian kernel more generally as

$$K(\mathbf{x}, \mathbf{y}) = e^{-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{y}\|_2^2} \quad (\sigma^2 > 0 \text{ to be tuned via CV})$$

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 2: Gaussian Kernel/Radial Basis Function Kernel



Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

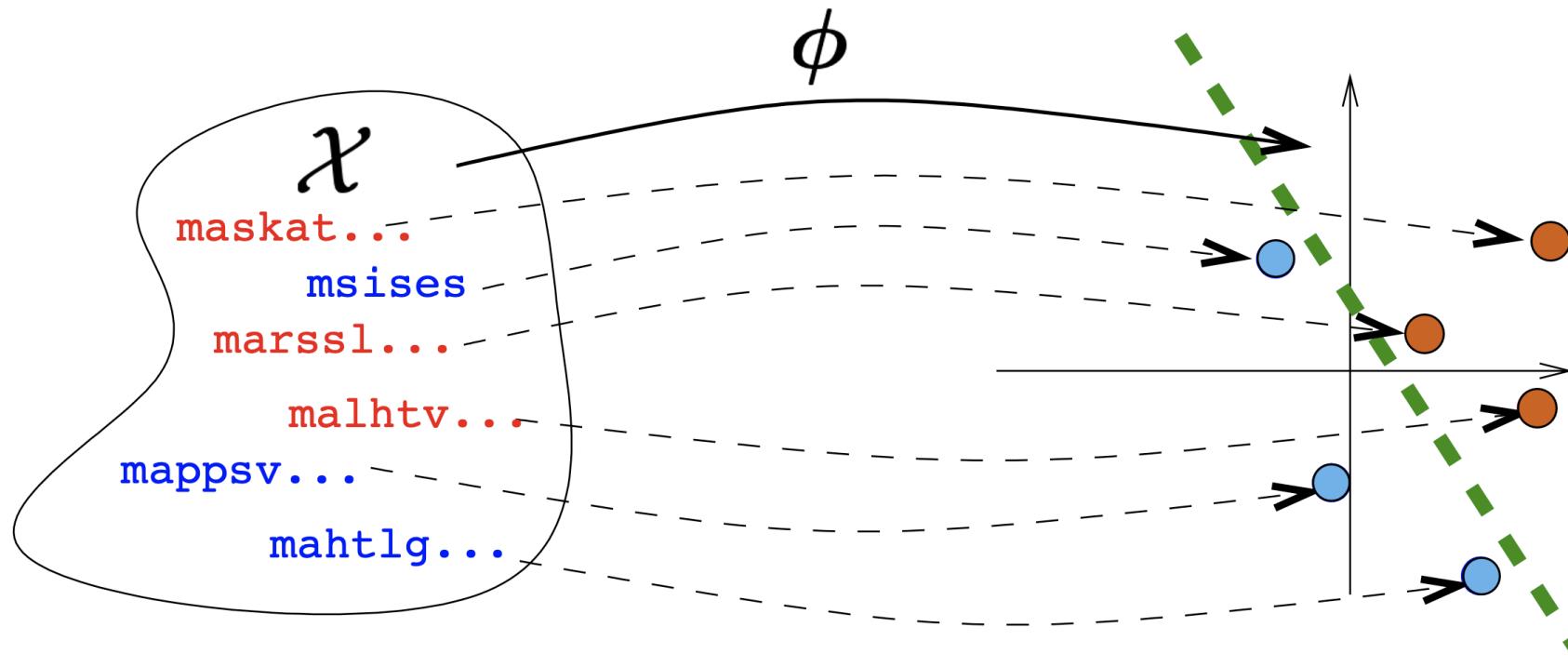
Example 3: String Kernels

- The examples given so far were all for kernels over Euclidean vector spaces.
In many learning tasks found in practice, the input space \mathcal{X} is **not** a vector space.
 - A few examples: classify protein sequences, images, graphs, parse trees, finite automata, or other discrete structures which may not be directly given as vectors.
- The kernel tricks provide a method for extending algorithms such as SVMs originally designed for a vectorial space to the classification of such objects.
- Here we will focus on the specific case of *sequence kernels*, that is, kernels for *sequences* or *strings*.
 - Kernels can be defined for other discrete structures in somewhat similar ways.
 - *Sequence kernels* are particularly relevant to learning algorithms applied to computational biology or natural language processing, which are both important applications.

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 3: *String Kernels*



Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 3: String Kernels

- Kernels, i.e. similarity measures, for sequences, can be based on the idea of declaring two strings, e.g., two documents or two biosequences, as similar when they share **common substrings** or **subsequences**.
- One example could be the kernel between two sequences defined by the *sum of the product of the counts of their common substrings*.
 - But which substrings should be used in that definition?
 - Most likely, we would need some flexibility in the definition of the matching substrings.
 - For computational biology applications, for example, the match could be imperfect. Thus, we may need to consider some number of mismatches, possibly gaps, or wildcards.
 - More generally, we might need to allow various substitutions and might wish to assign different weights to common substrings to emphasize some matching substrings and deemphasize others.
- Hence, as it is clear, there are many different possibilities and frameworks for defining such kernels.
- Here we will consider **spectrum kernel** and **substring kernel**. For a more general approach based on *rational kernel*, see Section 6.5 of [Mohri et al. \(2018\)](#).

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 3: String Kernels

General Approach

- Remap the string-space \mathcal{X} into *fixed-length strings*; that is, for each string $\boldsymbol{x} \in \mathcal{X}$ define the *feature map*

$$\phi(\boldsymbol{x}) = [\phi_u(\boldsymbol{x})]_{u \in \mathcal{A}^k} \quad \text{with} \quad \mathcal{A}^k = \{\text{strings of length } k \text{ from a given alphabet } \mathcal{A}\}$$

where $\phi_u(\boldsymbol{x})$ can be:

- The number of occurrences of u in \boldsymbol{x} (without gaps) \rightsquigarrow **spectrum kernel** (Leslie et al., 2002)
- The number of occurrences of u in \boldsymbol{x} up to m mismatches (without gaps) \rightsquigarrow **mismatch kernel** (Leslie et al., 2004)
- The number of occurrences of u in \boldsymbol{x} allowing gaps, with a weight decaying exponentially with the number of gaps \rightsquigarrow **all-substrings kernel** (Lohdi et al., 2002)

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 3a: *Spectrum Kernels* = {count how many (contiguous) length- k substrings \mathbf{x} and \mathbf{y} share}

- For example, the **3-spectra** (with no gaps) of $\mathbf{x} = \text{STATISTICS}$ and $\mathbf{y} = \text{COMPUTATION}$ are:
 $(\text{STA}, \text{TAT}, \text{ATI}, \text{TIS}, \text{IST}, \text{STI}, \text{TIC}, \text{ICS})$ and $(\text{COM}, \text{OMP}, \text{MPU}, \text{PUT}, \text{UTA}, \text{TAT}, \text{ATI}, \text{TIO}, \text{ION})$

They have in common the substrings **TAT** and **ATI** $\rightsquigarrow K(\mathbf{x}, \mathbf{y}) = 2$

- More in general, the **k -spectrum kernel** is then:

$$K(\mathbf{x}, \mathbf{y}) = \sum_{u \in \mathcal{A}^k} \phi_u(\mathbf{x}) \phi_u(\mathbf{y}) \quad \text{with} \quad \phi_u(\mathbf{x}) = \{\text{Number of occurrences of } u \text{ in } \mathbf{x}\}$$

- The computation of the kernel is formally a sum over $|\text{card}(\mathcal{A})|^k$ terms, but at most $\text{card}(\mathbf{x}) - k + 1$ terms are non-zero in $\phi(\mathbf{x}) \rightsquigarrow$ computation in $O(\text{card}(\mathbf{x}) + \text{card}(\mathbf{y}))$ with strings pre-indexation.
 - Fast classification of a sequence \mathbf{x} in $O(\text{card}(\mathbf{x}))$
 - Work with any string (natural language, time series, etc) + **fast** and **scalable**

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 3b: All-Substrings Kernels

- For $1 \leq k \leq m \in \mathbb{N}$, denote by $\mathcal{I}(k, m)$ the set of **sequences of indices** $\mathbf{i} = (i_1, \dots, i_k)$, with $1 \leq i_1 < i_2 < \dots < i_k \leq m$
- For a length m string $\mathbf{x} = x_1 x_2 \cdots x_m \in \mathcal{X}$, and a sequence of indices $\mathbf{i} \in \mathcal{I}(k, m)$, define a *substring* as

$$\mathbf{x}(\mathbf{i}) = x_{i_1} x_{i_2} \cdots x_{i_k} \quad \text{and} \quad \ell(\mathbf{i}) = \{\text{length of } \mathbf{x}(\mathbf{i})\} = i_k - i_1 + 1$$

- For example, if $\mathbf{x} = \text{ABRACADABRA}$, then

$$\mathbf{i} = (3, 4, 7, 8, 10) \quad \rightsquigarrow \quad \mathbf{x}(\mathbf{i}) = \text{RADAR} \quad \rightsquigarrow \quad \ell(\mathbf{i}) = 10 - 3 + 1 = 8$$

- Let $k \in \mathbb{N}$ and $\lambda > 0$ be tuning parameter to be chosen via CV \rightsquigarrow let $\phi_u : \mathcal{X} \mapsto \mathbb{R}$ be defined as

$$\phi_u(\mathbf{x}) = \sum_{\substack{\mathbf{i} \in \mathcal{I}(k, \text{card}(\mathbf{x})) \\ \text{s.t. } \mathbf{x}(\mathbf{i})=u}} \lambda^{\ell(\mathbf{i})} \quad \rightsquigarrow \text{substring kernel} \quad K(\mathbf{x}, \mathbf{y}) = \sum_{u \in \mathcal{A}^k} \phi_u(\mathbf{x}) \phi_u(\mathbf{y})$$

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 3b: All-Substrings Kernels

u	ca	ct	at	ba	bt	cr	ar	br
$\phi_u(\text{cat})$	λ^2	λ^3	λ^2	0	0	0	0	0
$\phi_u(\text{car})$	λ^2	0	0	0	0	λ^3	λ^2	0
$\phi_u(\text{bat})$	0	0	λ^2	λ^2	λ^3	0	0	0
$\phi_u(\text{bar})$	0	0	0	λ^2	0	0	λ^2	λ^3

- For example:

- $K(\text{cat}, \text{cat}) = K(\text{car}, \text{car}) = 2\lambda^4 + \lambda^6$
- $K(\text{cat}, \text{car}) = \lambda^4$
- $K(\text{cat}, \text{bar}) = 0$

Computations

$$K(\mathbf{x}, \mathbf{y}) = \sum_{u \in \mathcal{A}^k} \phi_u(\mathbf{x}) \phi_u(\mathbf{y}) = \sum_{u \in \mathcal{A}^k} \sum_{\mathbf{i}: \mathbf{x}(\mathbf{i})=u} \sum_{\mathbf{i}' : \mathbf{y}(\mathbf{i}')=u} \lambda^{\ell(\mathbf{i})+\ell(\mathbf{i}')}$$

Enumerating the substrings is of order $O(\text{card}(\mathbf{x})^k)$ \rightsquigarrow **too slow** \rightsquigarrow but we can speed it up (see Ch. 11 of Shawe-Taylor and Cristianini, 2004)

Toward Nonlinearity...and Beyond!

B) Kernel-trick \longleftrightarrow Implicit transformation

Example 4: Kernel for Text

- Probably the most important data type after vectors and sequences/strings is **free text**.
- Kernel methods provide a natural framework for pattern analysis in text \rightsquigarrow here similarity measures can often be obtained by borrowing methods from different disciplines
- Well-known techniques from *Information Retrieval* (IR), such as the rich class of **vector space models**, (VSM) can naturally be interpreted as kernel methods.
- An important property of the vector space representation is that the **primal-vs-dual** framework we have developed for SVM has an interesting counterpart in the interplay between **word-based-vs-document-based** representations.
- Other kernel constructions like *string-matching* (see above) can be applied to free texts, but clearly they are **not** specific to *natural language text*.
- Here we will consider the most basic VSMs only. For mode examples and comments see Ch. 10 of **Shawe-Taylor and Cristianini, 2004**)

Toward Nonlinearity...and Beyond!

B) Kernel-trick \leftrightarrow Implicit transformation

Example 4: *Kernel for Text* \rightsquigarrow *Vector Space Models* (VSM) \rightsquigarrow *Bag-of-Words* (BoW)

- BoW is the simplest possible version of VSM \rightsquigarrow a *document* is represented by the *words* it contains (and their frequency) with the ordering and punctuation ignored
 - *Words/Terms* (tr): sequence of letters from an alphabet \mathcal{A} separated by punctuation or spaces
 - *Document* (doc): arbitrary collection of words
 - *Corpus* (\mathcal{C}_n): set of n documents we are processing
 - *Dictionary* (\mathcal{D}_p): a permanent set of p terms or, more often, the set of p words in the *corpus*
- A BoW can be represented as a vector: **each component** is associated with **one term** from the dictionary

$$\phi(\text{doc}) = [\text{tf}(\text{tr}_1, \text{doc}) \quad \text{tf}(\text{tr}_2, \text{doc}) \quad \dots \quad \text{tf}(\text{tr}_p, \text{doc})]^T \in \mathbb{R}^p$$

where $\text{tf}(\text{tr}_j, \text{doc}) = \{\text{frequency of the } j^{\text{th}} \text{ term in the document } \text{doc}\} \rightsquigarrow \text{sparse but very high-dimensional embedding!}$

Toward Nonlinearity...and Beyond!

B) Kernel-trick \leftrightarrow Implicit transformation

Example 4: Kernel for Text \rightsquigarrow Vector Space Models (VSM) \rightsquigarrow Bag-of-Words (BoW)

- Document-Term Matrix (DTM) \rightsquigarrow rows = docs = obs / cols = words = variables ($n \times p$)

$$\mathbb{X} = \begin{bmatrix} \phi(\text{doc}_1)^T \\ \phi(\text{doc}_2)^T \\ \vdots \\ \phi(\text{doc}_n)^T \end{bmatrix} = \begin{bmatrix} \text{tf}(\text{tr}_1, \text{doc}_1) & \cdots & \text{tf}(\text{tr}_p, \text{doc}_1) \\ \text{tf}(\text{tr}_1, \text{doc}_2) & \cdots & \text{tf}(\text{tr}_p, \text{doc}_2) \\ \vdots & \ddots & \vdots \\ \text{tf}(\text{tr}_1, \text{doc}_n) & \cdots & \text{tf}(\text{tr}_p, \text{doc}_n) \end{bmatrix}$$

$$\begin{array}{ccccccccc} \mathbb{X}^T & & \mathbb{X}^T \mathbb{X} & & \mathbb{X} \mathbb{X}^T & \rightsquigarrow & \mathbb{K} = \mathbb{X} \mathbb{X}^T & \rightsquigarrow & K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathbb{R}^p} = \sum_{j=1}^p \text{tf}(\text{tr}_j, \mathbf{x}) \cdot \text{tf}(\text{tr}_j, \mathbf{y}) \\ \text{Term-Doc} & & \text{Term-Term} & & \text{Doc-Doc} & & \text{kernel-matrix} & & \\ \text{TDM: } (p \times n) & & \text{TTM: } (p \times p) & & \text{DDM: } (n \times n) & & \text{(Gramian)} & & \end{array}$$

- **Interpreting Duality:** the *dual* corresponds to a doc view, while the *primal* provides a term view

$\text{doc} = \{\text{counts of the terms that appear in it}\}$ vs $\text{tr} = \{\text{counts of the docs in which it appears}\}$

- **Remark:** often $n \ll p \rightsquigarrow$ dual/docs computationally better vs primal/terms better for interpretation

Toward Nonlinearity...and Beyond!

B) Kernel-trick \leftrightarrow Implicit transformation

Example 4: Kernel for Text \rightsquigarrow Vector Space Models (VSM) \rightsquigarrow Bag-of-Words (BoW)

Improvements and Limitations

The refinements below can be performed *sequentially* and thought of as a series of successive *embeddings* \rightsquigarrow we end up with a single but *composite* feature map that incorporates different aspects of **domain knowledge**.

- *Normalization*: the longer the **doc**, the more words it contains \rightsquigarrow the larger $\|\phi(\text{doc})\|$ is \rightsquigarrow normalize!

$$\widetilde{K}(\mathbf{x}, \mathbf{y}) = \left\langle \frac{\phi(\mathbf{x})}{\|\phi(\mathbf{x})\|}, \frac{\phi(\mathbf{y})}{\|\phi(\mathbf{y})\|} \right\rangle_{\mathbb{R}^p} = \frac{K(\mathbf{x}, \mathbf{y})}{\sqrt{K(\mathbf{x}, \mathbf{x})} \cdot \sqrt{K(\mathbf{y}, \mathbf{y})}} \quad (\text{typically the first or the last transformation})$$

- *Nonlinear Embeddings*: starting from the baseline kernel $K(\mathbf{x}, \mathbf{y})$ corresponding to the (normalised) BoW representation, we can always go **nonlinear** using standard kernel construction. For example

$$\widetilde{K}(\mathbf{x}, \mathbf{y}) = (c + K(\mathbf{x}, \mathbf{y}))^d \quad (\text{polynomial-kernel over BoW representation})$$

uses all k of words for $0 \leq k \leq d$ as features.

Toward Nonlinearity...and Beyond!

B) Kernel-trick \leftrightarrow Implicit transformation

Example 4: Kernel for Text \rightsquigarrow Vector Space Models (VSM) \rightsquigarrow Bag-of-Words (BoW)

Improvements and Limitations

- *Term weighting*: not all words have the same importance in a task (e.g. guess the *topic* of a doc) \rightsquigarrow stop-words are removed before the analysis starts \rightsquigarrow equivalent to assign weight 0 to these coordinates in ϕ . There are *relative* (to a topic or task) or *absolute* measures of importance \rightsquigarrow the latter do **not** use label/response info \rightsquigarrow can be pre-computed on external-larger-unlabeled corpora.
 - $\text{idf} = \{\text{inverse doc frequency}\}$ is an *absolute* measure \rightsquigarrow for any word $\text{tr} \in \mathcal{D}_p$, let

$$\text{df}(\text{tr}) = \{\text{number of docs in } \mathcal{C}_n \text{ containing } \text{tr}\} \rightsquigarrow \text{idf}(\text{tr}) = \ln\left(\frac{n}{\text{df}(\text{tr})}\right)$$

$$\mathbb{R} = \begin{bmatrix} \text{idf}(\text{tr}_1) & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \text{idf}(\text{tr}_p) \end{bmatrix} \rightsquigarrow \widetilde{K}(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) [\mathbb{R} \mathbb{R}^T] \phi(\mathbf{y})^T = \sum_{j=1}^p \text{idf}(\text{tr}_j)^2 \text{tf}(\text{tr}_j, \mathbf{x}) \text{tf}(\text{tr}_j, \mathbf{y})$$

(tf-idf representation)

log-scale \rightsquigarrow no weight is too large relative to the weight of a term that occurs in roughly half the docs

Toward Nonlinearity...and Beyond!

B) Kernel-trick \leftrightarrow Implicit transformation

Example 4: *Kernel for Text* \rightsquigarrow *Vector Space Models* (VSM) \rightsquigarrow *Bag-of-Words* (BoW)

Improvements and Limitations

- *Term weighting*: not all words have the same importance in a task (e.g. guess the *topic* of a doc) \rightsquigarrow stop-words are removed before the analysis starts \rightsquigarrow equivalent to assign *weight* 0 to these coordinates in ϕ
- *Semantic issues*: **no** semantic relation between words is extraced \rightsquigarrow homonyms/synonyms = problems tf-idf down-weights irrelevant terms + highlights potentially discriminative ones \rightsquigarrow **no** semantic!
 - *Stemming*: removing **inflections** from words ensures that different forms of the same words are treated as equivalent terms \rightsquigarrow basic type of semantic similarity implemented as part of the **tokeniser**.
 - **Solution**: explicitly introduce a **semantic similarity weight** \rightsquigarrow a $(p \times p)$ *proximity matrix* \mathbb{P} such that $\mathbb{P}[i, j] > 0 \Leftrightarrow \mathbf{tr}_i$ is *semantically* related to $\mathbf{tr}_j \rightsquigarrow \widetilde{K}(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})[\mathbb{P}\mathbb{P}^T]\phi(\mathbf{y})^T$.
 - Many methods are available to design proximity matrices by analysing **co-occurrence** information.
For example: *generalized vector space models* (GVSM), *latent semantic*, *semantic diffusion kernels*, etc.
 - *Remember*: NN-based embeddings for text mentioned earlier \rightsquigarrow **BERT** and **similia**, **FastText**, etc.

Happy ending...

Summing up

- SVMs can be *really* good classifiers/predictors; from about 1995--2012 they won all the contests and they might come around again...
- Following Vapnik's reductionism, **regularization** and **surrogate loss**, two pillars of modern learning theory, flow naturally from purely *geometric* arguments.
- The concept of margin plays a crucial rule and, typically, classifiers with large margin generalize well
 - Asking for a large margin is an extra constraint that limits the model complexity and *geometrically* code for a L^2 -penalty
 - Maximizing the margin is computationally efficient and leads to sparse solutions (support vectors)
 - If we can't get perfect classification we allow ourselves some **slack**
- Nonlinearity can be nicely and efficiently brought into the picture by exploiting the *kernel-trick*.
 - Learning task where the input space \mathcal{X} is **not** a vector space can also be treated effortlessly and in a unified manner.

Going forward

...broaden our understanding of (Mercer) kernel methods in the framework of RKHS...

Reproducing Kernel Hilbert Spaces

References

