

CSCI 567: Machine Learning

Vatsal Sharan
Fall 2022

Lecture 4, Sep 15

Administrivia

- HW2 will be released tonight, due in about 2 weeks.
- We will post some practice problems for the quiz by early next week.

Recap

Ensuring generalization

Theorem. Let \mathcal{F} be a function class with size $|\mathcal{F}|$. Let $y = f^*(\mathbf{x})$ for some $f^* \in \mathcal{F}$. Suppose we get a training set $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ of size n with each datapoint drawn i.i.d. from the data distribution D . Let

$$f_S^{ERM} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i).$$

For any constants $\epsilon, \delta \in (0, 1)$, if $n \geq \frac{\ln(|\mathcal{F}|/\delta)}{\epsilon}$, then with probability $(1 - \delta)$ over $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, $R(f_S^{ERM}) < \epsilon$.

A useful rule of thumb: to guarantee generalization, make sure that your training data set size n is at least linear in the number d of free parameters in the function that you're trying to learn.

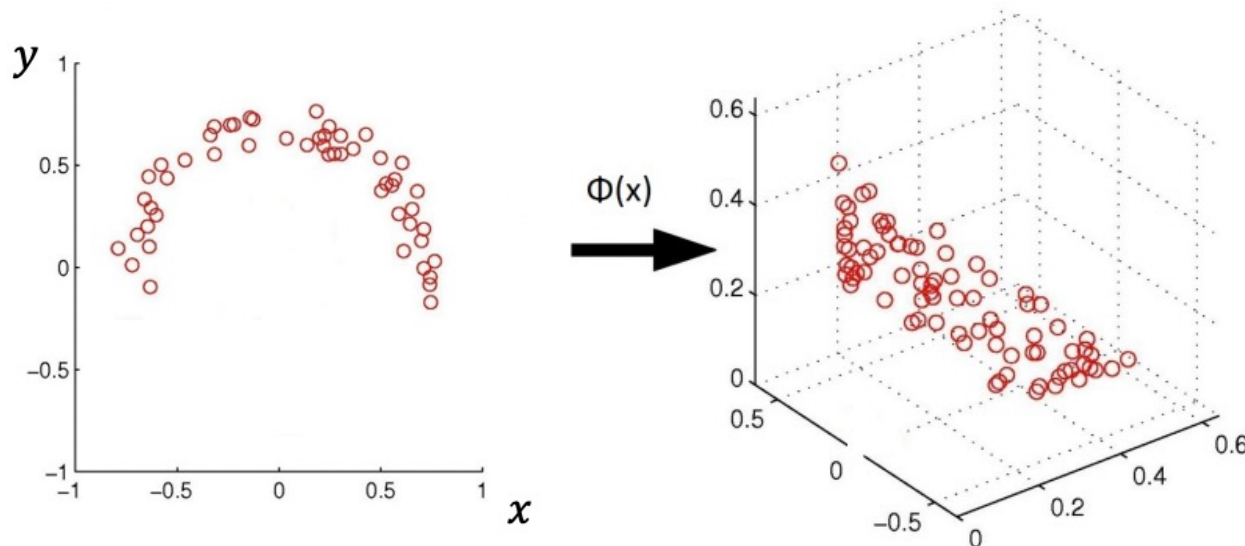
Beyond linear models: nonlinearly transformed features

1. Use a nonlinear mapping

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d \rightarrow \mathbf{z} \in \mathbb{R}^M$$

to transform the data to a more complicated feature space

2. Then apply linear regression (hope: linear model is a better fit for the new feature space).



Polynomial basis functions

Polynomial basis functions for $d = 1$

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

Learning a linear model in the new space

= learning an *M -degree polynomial model* in the original space

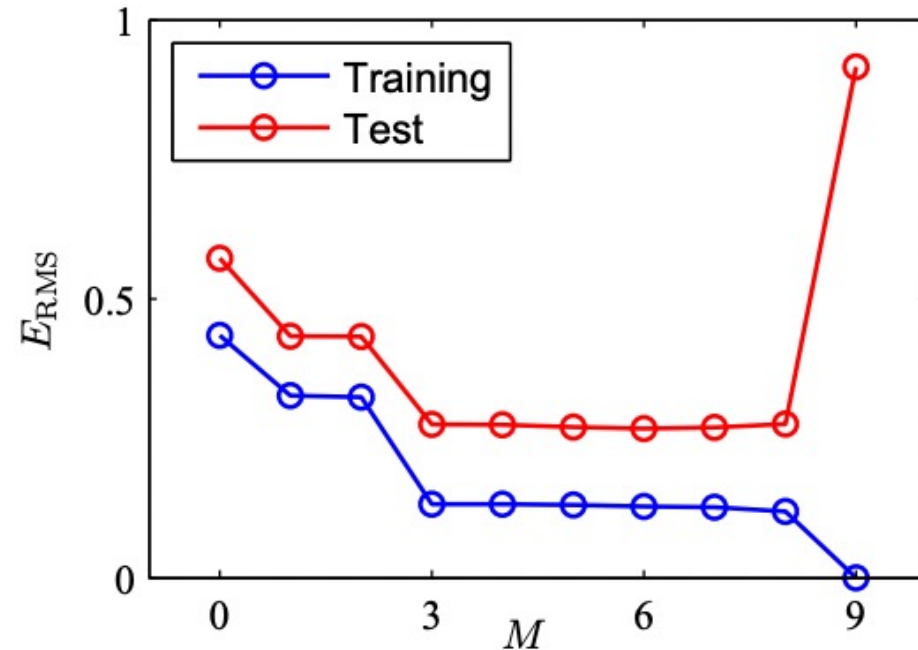
Underfitting and overfitting

$M \leq 2$ is *underfitting* the data

- large training error
- large test error

$M \geq 9$ is *overfitting* the data

- small training error
- **large test error**



More complicated models \Rightarrow larger gap between training and test error

How to prevent overfitting?

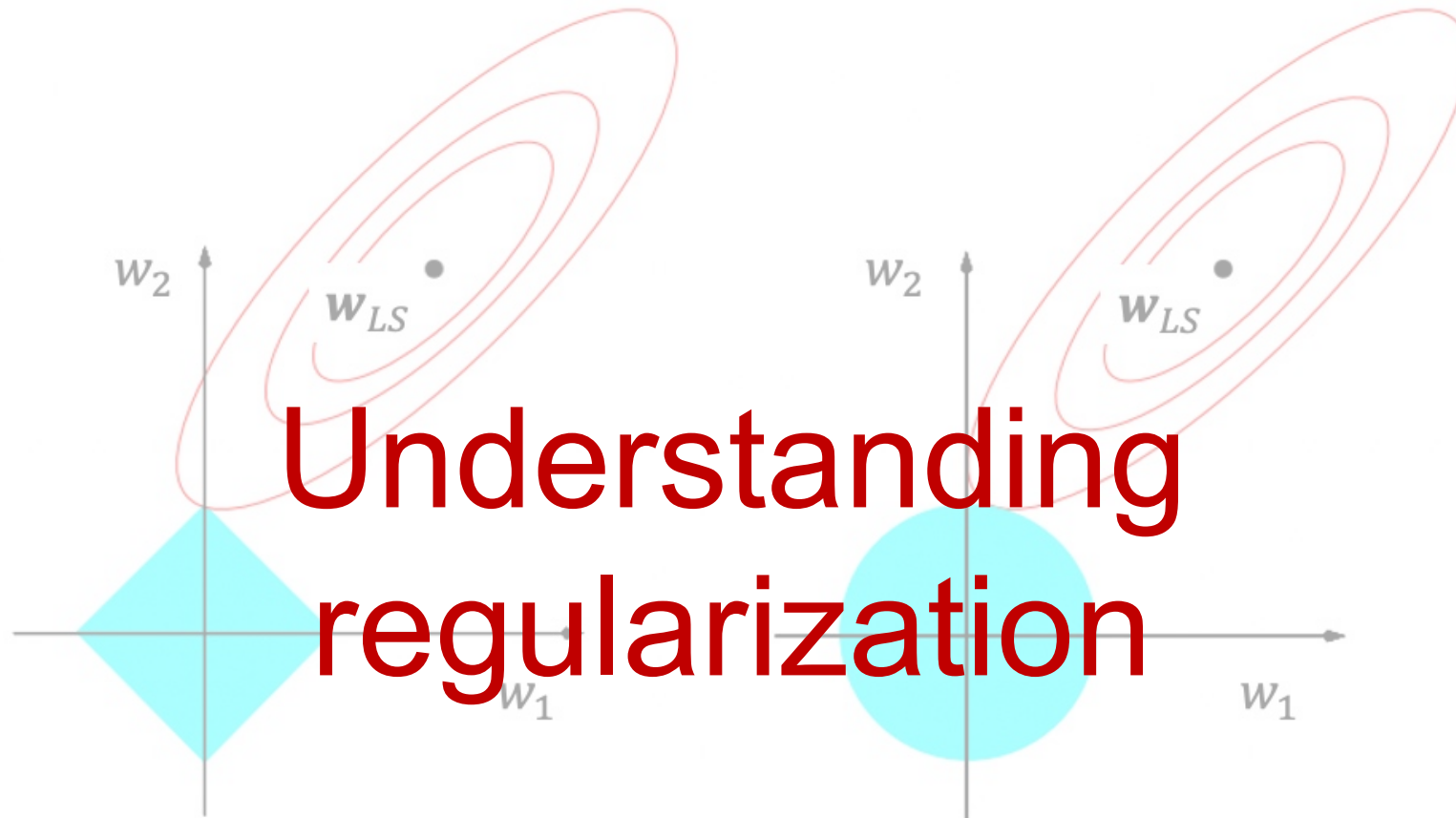
Preventing overfitting: **Regularization**

Regularized linear regression: new objective

$$G(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda\psi(\mathbf{w})$$

Goal: find $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} G(\mathbf{w})$

- $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^+$ is the *regularizer*
 - measure how complex the model \mathbf{w} is, penalize complex models
 - common choices: $\|\mathbf{w}\|_2^2$, $\|\mathbf{w}\|_1$, etc.
- $\lambda > 0$ is the *regularization coefficient*
 - $\lambda = 0$, no regularization
 - $\lambda \rightarrow +\infty$, $\mathbf{w} \rightarrow \operatorname{argmin}_{\mathbf{w}} \psi(\mathbf{w})$
 - i.e. control **trade-off** between training error and complexity



ℓ_2 regularization: penalizing large weights

ℓ_2 regularization, $\psi(\mathbf{w}) = \|\mathbf{w}\|_2^2$:

$$G(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\nabla G(\mathbf{w}) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) + 2\lambda \mathbf{w} = 0$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Linear regression with ℓ_2 regularization is also known as **ridge regression**.

With a Bayesian viewpoint, corresponds to a Gaussian prior for \mathbf{w} .

Encouraging sparsity: ℓ_0 regularization

Sparsity of w : Number of non-zero coefficients in w . Same as $\|w\|_0$

Advantage:

- Sparse models are a natural inductive bias in many settings. In many applications we have numerous possible features, only some of which may have any relationship with the label.
- Sparse models may also be more **interpretable**. They could narrow down a small number of features which carry a lot of signal.
- Data required to learn sparse model maybe significantly less than to learn dense model.

We'll see more on the third point next.

ℓ_0 regularization: The good, the bad and the ugly

Choose $\psi(\boldsymbol{w}) = \|\boldsymbol{w}\|_0$.

$$G(\boldsymbol{w}) = \sum_{i=1}^n (\boldsymbol{w}^T \boldsymbol{x}_i - y_i)^2 + \lambda \|\boldsymbol{w}\|_0.$$

ℓ_0 regularization: The good, the bad and the ugly

ℓ_0 regularization: The good, the bad and the ugly

ℓ_0 regularization: The good, the bad and the ugly

ℓ_1 regularization as a proxy for ℓ_0 regularization

Choose $\psi(\boldsymbol{w}) = \|\boldsymbol{w}\|_1$.

$$G(\boldsymbol{w}) = \sum_{i=1}^n (\boldsymbol{w}^T \boldsymbol{x}_i - y_i)^2 + \lambda \|\boldsymbol{w}\|_1.$$

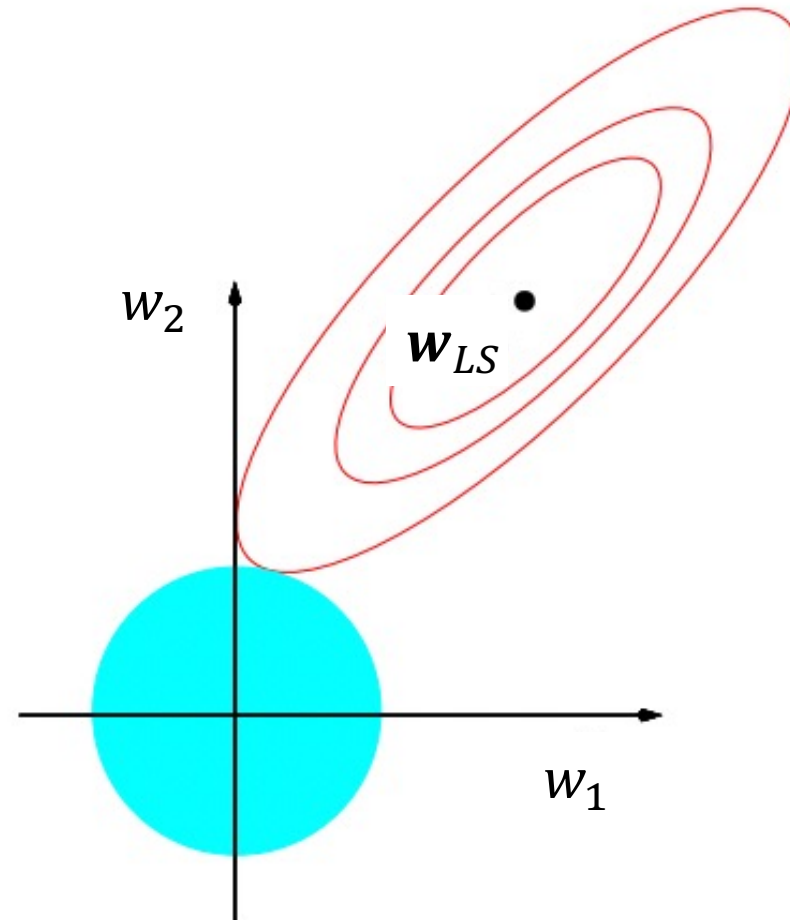
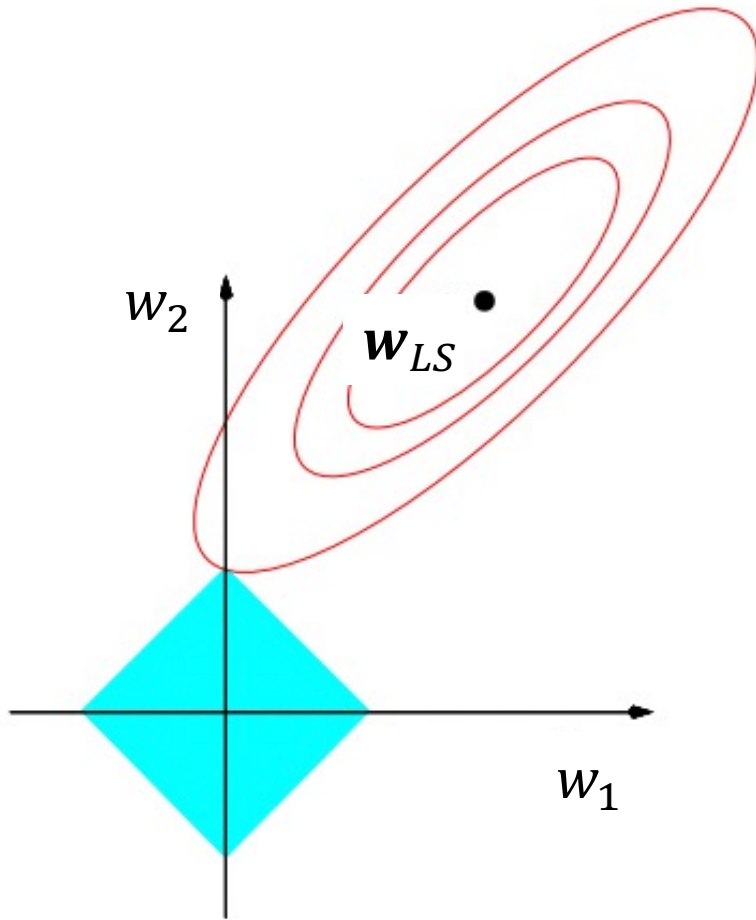
ℓ_1 regularization as a proxy for ℓ_0 regularization

Theorem. *Given n vectors $\{\mathbf{x}_i \in \mathbb{R}^d, i \in [n]\}$ drawn i.i.d. from $N(0, \mathbf{I})$, let $y_i = \mathbf{w}^{*T} \mathbf{x}_i$ for some \mathbf{w}^* with $\|\mathbf{w}^*\|_0 = s$. Then for some fixed constant $C > 0$, the minimizer of $G(\mathbf{w})$ with $\psi(\mathbf{w}) = \|\mathbf{w}\|_1$ will be \mathbf{w}^* as long as $n > C \cdot s \log d$ (with high probability over the randomness in the training datapoints \mathbf{x}_i).*

[similar result can also be proven under more general conditions].

Why does ℓ_1 regularization encourage sparse solutions?

Optimization problem: $\operatorname{argmin}_{\mathbf{w}} \operatorname{RSS}(\mathbf{w})$, subject to $\psi(\mathbf{w}) \leq \beta$



Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Let $\beta_j = \mathbf{X}_{(i)}^T \mathbf{y}$

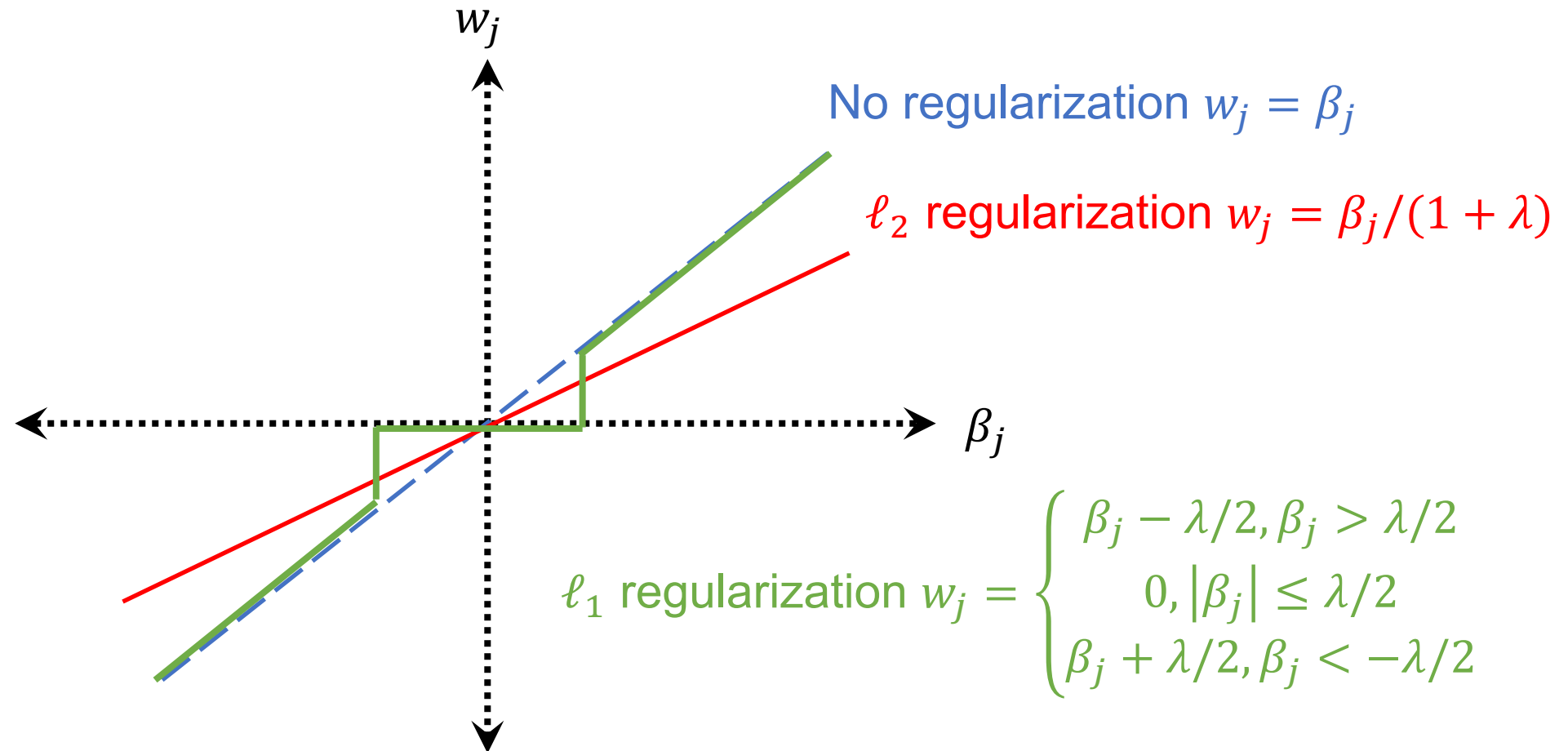
Using subgradients, we can show that for the ℓ_1 regularized case:

$$w_j = \begin{cases} \beta_j - \lambda/2, & \beta_j > \lambda/2 \\ 0, & |\beta_j| \leq \lambda/2 \\ \beta_j + \lambda/2, & \beta_j < -\lambda/2 \end{cases}$$

Diving deeper: ℓ_1 and ℓ_2 regularization for the “isotropic” case

Summary: Isotropic case ($\mathbf{X}^T \mathbf{X} = \mathbf{I}$).

Let $\beta_j = \mathbf{X}_{(i)}^T \mathbf{y}$



Implicit regularization

So far, we explicitly added a $\psi(\mathbf{w})$ term to our objective function to regularize.

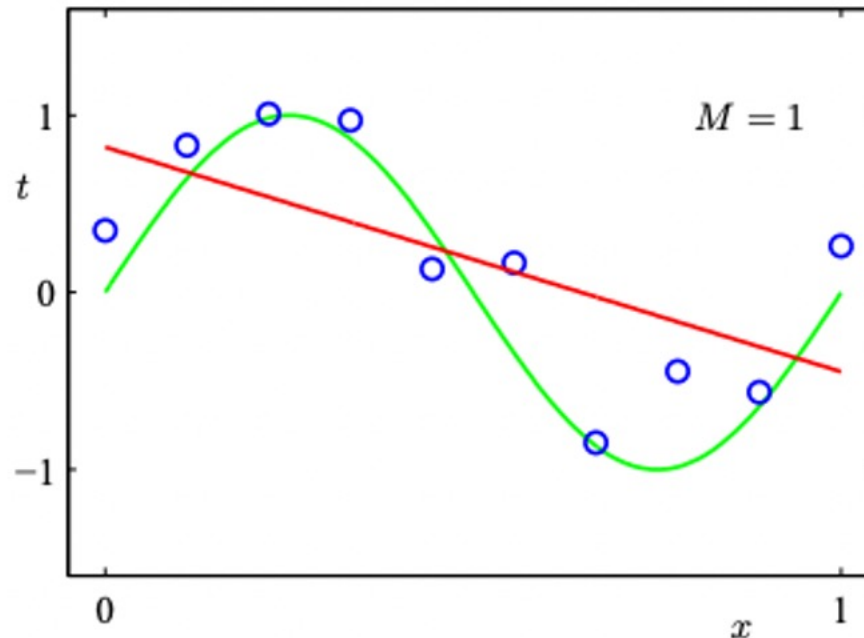
In many cases, the optimization algorithm we use can themselves act as regularizers, favoring some solutions over others.

Currently a very active area of research, you'll see more in the homework.

Bias-variance tradeoff

The phenomenon of underfitting and overfitting is often referred to as the *bias-variance tradeoff* in the literature.

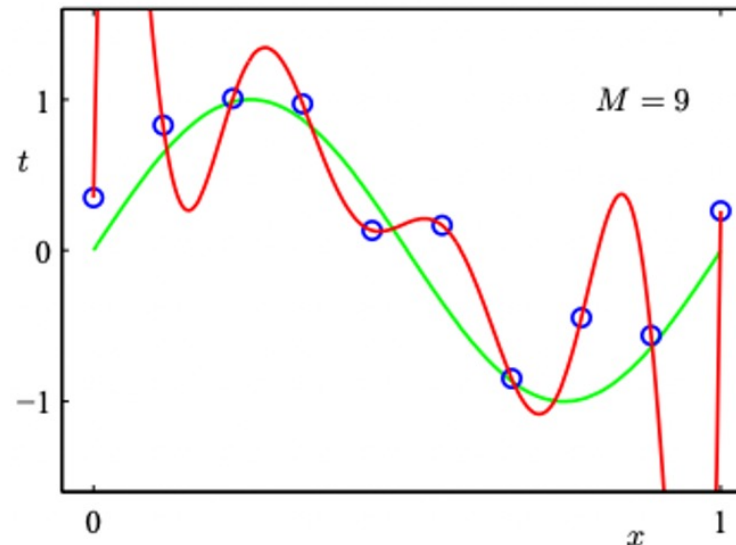
A model whose complexity is too *small* for the task will *underfit*. This is a model with a large bias because the model's accuracy will not improve even if we add a lot of training data.

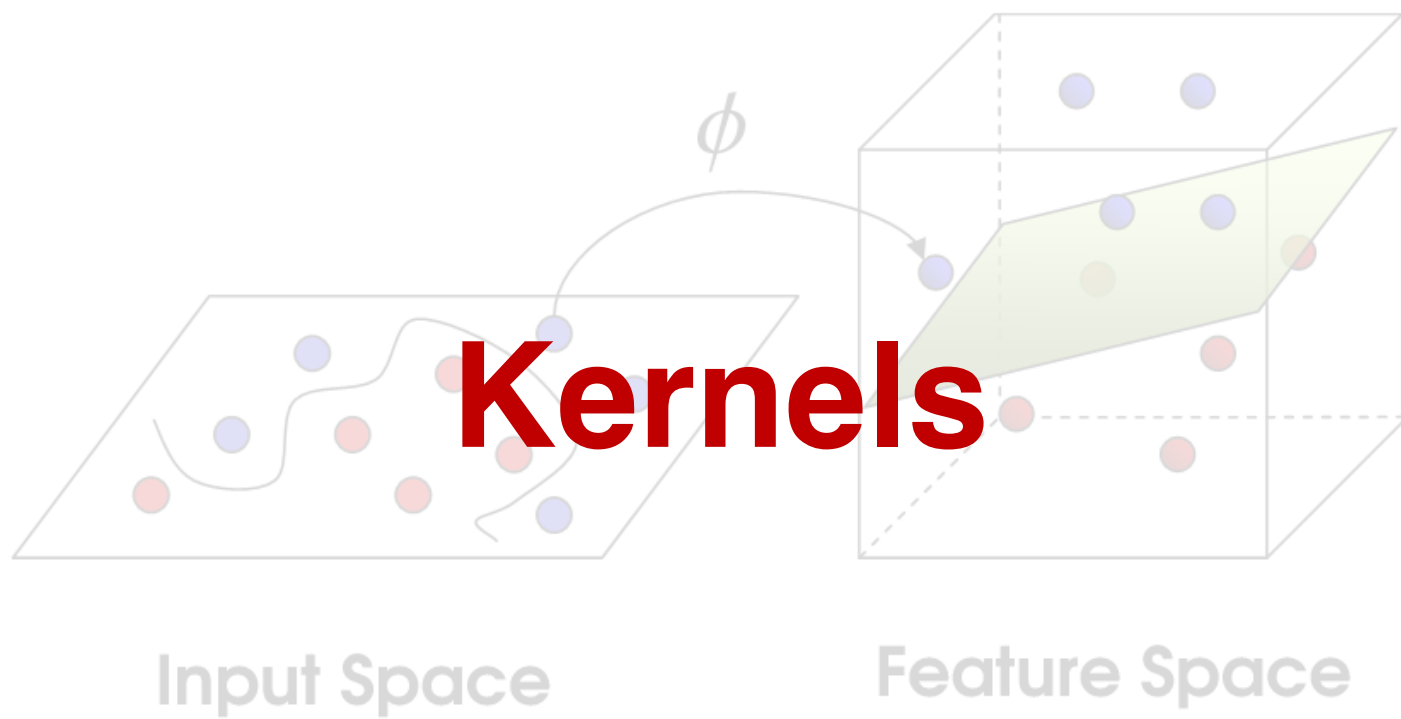


Bias-variance tradeoff

The phenomenon of underfitting and overfitting is often referred to as the ***bias-variance tradeoff*** in the literature.

A model whose complexity is too *large* for the amount of available training data will *overfit*. This is a model with high variance, because the model's predictions will vary a lot with the randomness in the training data (it can even fit any noise in the training data).





Motivation

Recall the nonlinear function map for linear regression:

1. Use a nonlinear mapping

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d \rightarrow \mathbf{z} \in \mathbb{R}^M$$

to transform the data to a more complicated feature space

2. Then apply linear regression (hope: linear model is a better fit for the new feature space).

Kernel methods give a way to choose and efficiently work with the nonlinear map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^M$ (for linear regression, and much more broadly).

Regularized least squares

Let's continue with regularized least squares with non-linear basis:

$$\begin{aligned}\mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmin}} F(\mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \left(\|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \right) \\ &= \left(\Phi^T \Phi + \lambda \mathbf{I} \right)^{-1} \Phi^T \mathbf{y}\end{aligned}$$

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_n)^T \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

This operates in space \mathbb{R}^M and M could be huge (and even infinite).

Regularized least squares solution: Another look

By setting the gradient of $F(\mathbf{w}) = \|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$ to be $\mathbf{0}$:

$$\Phi^T(\Phi\mathbf{w}^* - \mathbf{y}) + \lambda\mathbf{w}^* = \mathbf{0}$$

we know

$$\mathbf{w}^* = \frac{1}{\lambda}\Phi^T(\mathbf{y} - \Phi\mathbf{w}^*) = \Phi^T\boldsymbol{\alpha} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$$

Thus the least square solution is **a linear combination of features of the datapoints!**

This calculation does not show what $\boldsymbol{\alpha}$ should be, but ignore that for now.

Why is this helpful?

Assuming we know α , the prediction of w^* on a new example x is

$$w^{*\top} \phi(x) = \sum_{i=1}^n \alpha_i \phi(x_i)^\top \phi(x)$$

Therefore, *only inner products in the new feature space matter!*

Kernel methods are exactly about computing inner products *without explicitly computing ϕ* .

But we need to figure out what α is first!

Solving for α , Step 1: Kernel matrix

Plugging in $w = \Phi^T \alpha$ into $F(w)$ gives

$$\begin{aligned} H(\alpha) &= F(\Phi^T \alpha) \\ &= \|\Phi \Phi^T \alpha - \mathbf{y}\|_2^2 + \lambda \|\Phi^T \alpha\|_2^2 \\ &= \|\mathbf{K} \alpha - \mathbf{y}\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha \mathbf{K} \quad (\mathbf{K} = \Phi \Phi^T \in \mathbb{R}^{n \times n}) \end{aligned}$$

\mathbf{K} is called **Gram matrix** or **kernel matrix** where the (i, j) -th entry is

$$\mathbf{K}_{(i,j)} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Kernel matrix: Example

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

Gram/Kernel matrix

$$\begin{aligned} \mathbf{K} &= \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \phi(x_1)^T \phi(x_3) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \phi(x_2)^T \phi(x_3) \\ \phi(x_3)^T \phi(x_1) & \phi(x_3)^T \phi(x_2) & \phi(x_3)^T \phi(x_3) \end{pmatrix} \\ &= \begin{pmatrix} 4 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 4 \end{pmatrix} \end{aligned}$$

Kernel matrix vs Covariance matrix

	dimensions	entry (i, j)	property
$\Phi\Phi^T$			
$\Phi^T\Phi$			

Solving for α , Step 2: Minimize the dual

Minimize (the so-called *dual formulation*)

$$H(\alpha) = \|\mathbf{K}\alpha - \mathbf{y}\|_2^2 + \lambda\alpha^\top \mathbf{K}\alpha$$

Setting the derivative to 0 we have

$$\mathbf{0} = (\mathbf{K}^2 + \lambda\mathbf{K})\alpha - \mathbf{K}\mathbf{y} = \mathbf{K}((\mathbf{K} + \lambda\mathbf{I})\alpha - \mathbf{y})$$

Thus $\alpha = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$ is a minimizer and we obtain

$$\mathbf{w}^* = \Phi^\top \alpha = \Phi^\top (\mathbf{K} + \lambda\mathbf{I})^{-1} \mathbf{y}$$

Exercise: *are there other minimizers? and are there other \mathbf{w}^* 's?*

Comparing two solutions

Minimizing $F(w)$ gives $w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$

Minimizing $H(\alpha)$ gives $w^* = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y$

Note I has different dimensions in these two formulas.

Natural question: *are the two solutions the same or different?*

They have to be the same because $F(w)$ has a unique minimizer!

And they are:

$$\begin{aligned} & (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \\ &= (\Phi^T \Phi + \lambda I)^{-1} \Phi^T (\Phi \Phi^T + \lambda I) (\Phi \Phi^T + \lambda I)^{-1} y \\ &= (\Phi^T \Phi + \lambda I)^{-1} (\Phi^T \Phi \Phi^T + \lambda \Phi^T) (\Phi \Phi^T + \lambda I)^{-1} y \\ &= (\Phi^T \Phi + \lambda I)^{-1} (\Phi^T \Phi + \lambda I) \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y \\ &= \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y \end{aligned}$$

The kernel trick

If the solutions are the same, then what is the difference?

First, computing $(\Phi\Phi^T + \lambda\mathbf{I})^{-1}$ can be more efficient than computing $(\Phi^T\Phi + \lambda\mathbf{I})^{-1}$ when $n \leq M$.

More importantly, computing $\alpha = (\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$ also *only requires computing inner products in the new feature space!*

Now we can conclude that the exact form of $\phi(\cdot)$ is not essential; *all we need to do is know the inner products $\phi(\mathbf{x})^T\phi(\mathbf{x}')$.*

For some ϕ it is indeed possible to compute $\phi(\mathbf{x})^T\phi(\mathbf{x}')$ without computing/knowing ϕ . This is the *kernel trick*.

The kernel trick: Example 1

Consider the following polynomial basis $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$:

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

What is the inner product between $\phi(\mathbf{x})$ and $\phi(\mathbf{x}')$?

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 \end{aligned}$$

Therefore, *the inner product in the new space is simply a function of the inner product in the original space.*

The kernel trick: Example 2

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{2d}$ is parameterized by θ :

$$\phi_{\theta}(\mathbf{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_m) \\ \sin(\theta x_m) \end{pmatrix}$$

What is the inner product between $\phi_{\theta}(\mathbf{x})$ and $\phi_{\theta}(\mathbf{x}')$?

$$\begin{aligned} \phi_{\theta}(\mathbf{x})^T \phi_{\theta}(\mathbf{x}') &= \sum_{m=1}^d \cos(\theta x_m) \cos(\theta x'_m) + \sin(\theta x_m) \sin(\theta x'_m) \\ &= \sum_{m=1}^d \cos(\theta(x_m - x'_m)) \quad \text{(trigonometric identity)} \end{aligned}$$

Once again, *the inner product in the new space is a simple function of the features in the original space.*

The kernel trick: Example 3

Based on ϕ_θ , define $\phi_L : \mathbb{R}^d \rightarrow \mathbb{R}^{2d(L+1)}$ for some integer L :

$$\phi_L(\mathbf{x}) = \begin{pmatrix} \phi_0(\mathbf{x}) \\ \phi_{\frac{2\pi}{L}}(\mathbf{x}) \\ \phi_{2\frac{2\pi}{L}}(\mathbf{x}) \\ \vdots \\ \phi_{L\frac{2\pi}{L}}(\mathbf{x}) \end{pmatrix}$$

What is the inner product between $\phi_L(\mathbf{x})$ and $\phi_L(\mathbf{x}')$?

$$\begin{aligned} \phi_L(\mathbf{x})^\top \phi_L(\mathbf{x}') &= \sum_{\ell=0}^L \phi_{\frac{2\pi\ell}{L}}(\mathbf{x})^\top \phi_{\frac{2\pi\ell}{L}}(\mathbf{x}') \\ &= \sum_{\ell=0}^L \sum_{m=1}^d \cos\left(\frac{2\pi\ell}{L}(x_m - x'_m)\right) \end{aligned}$$

The kernel trick: Example 4

When $L \rightarrow \infty$, even if we cannot compute $\phi(x)$ (since it's a vector of *infinite dimension*), we can still compute inner product:

$$\begin{aligned}\phi_{\infty}(\mathbf{x})^T \phi_{\infty}(\mathbf{x}') &= \int_0^{2\pi} \sum_{m=1}^d \cos(\theta(x_m - x'_m)) d\theta \\ &= \sum_{m=1}^d \frac{\sin(2\pi(x_m - x'_m))}{x_m - x'_m}\end{aligned}$$

Again, a simple function of the original features.

Note that when using this mapping in linear regression, we are *learning a weight w^* with infinite dimension!*

Kernel functions

Definition: a function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is called a *kernel function* if there exists a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^M$ so that for any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$,

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

Examples we have seen

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^2$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_{m=1}^d \frac{\sin(2\pi(x_m - x'_m))}{x_m - x'_m}$$

Using kernel functions

Choosing a nonlinear basis ϕ becomes equivalent to choosing a kernel function.

As long as computing the kernel function is more efficient, we should apply the kernel trick.

Gram/kernel matrix becomes:

$$\mathbf{K} = \Phi \Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

In fact, k is a kernel if and only if \mathbf{K} is positive semidefinite for *any n and any $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$* (**Mercer theorem**).

- useful for proving that a function is not a kernel

Examples which are not kernels

Function

$$k(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2^2$$

is *not a kernel*, why?

If it is a kernel, the kernel matrix for two data points \mathbf{x}_1 and \mathbf{x}_2 :

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

must be positive semidefinite, *but is it?*

Properties of kernels

For any function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})f(\mathbf{x}')$ is a kernel.

If $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are kernels, then the following are also kernels:

- **conical combination**: $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$ if $\alpha, \beta \geq 0$
- **product**: $k_1(\cdot, \cdot)k_2(\cdot, \cdot)$
- **exponential**: $e^{k(\cdot, \cdot)}$
- ...

Verify using the definition of kernel!

Popular kernels

Polynomial kernel

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^T \boldsymbol{x}' + c)^M$$

for $c \geq 0$ and M is a positive integer.

What is the corresponding ϕ ?

Popular kernels

Gaussian kernel or Radial basis function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2} \right) \quad \text{for some } \sigma > 0.$$

What is the corresponding ϕ ?

Popular kernels

Gaussian kernel or Radial basis function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2} \right) \quad \text{for some } \sigma > 0.$$

What is the corresponding ϕ ?

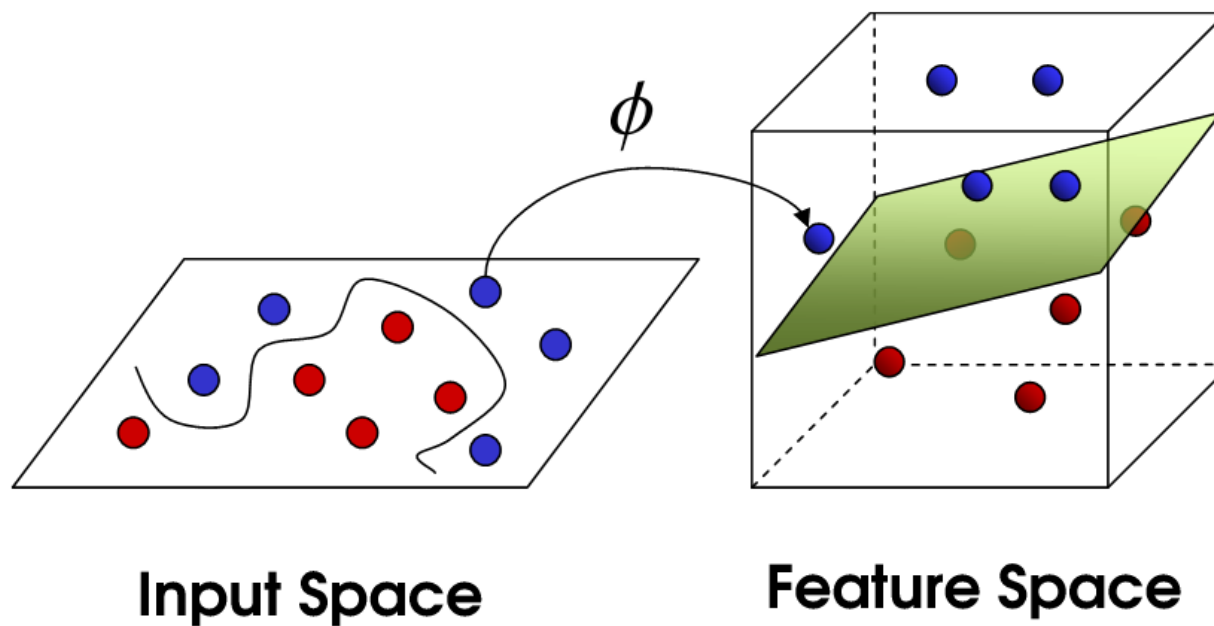
Prediction with kernels

As long as $\mathbf{w}^* = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$, prediction on a new example \mathbf{x} becomes

$$\mathbf{w}^{*\top} \phi(\mathbf{x}) = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}).$$

This is known as a **non-parametric method**. Informally speaking, this means that there is no fixed set of parameters that the model is trying to learn (remember \mathbf{w}^* could be infinite). Nearest-neighbors is another non-parametric method we have seen.

Classification with kernels



Similar ideas extend to the classification case, and we can predict using $\text{sign}(\mathbf{w}^T \phi)$.
Data may become linearly separable in the feature space!



A diagram illustrating a Support Vector Machine (SVM) classification. It shows two classes of data points: blue circles and red circles. A solid gray line represents the decision boundary. Two dashed gray lines parallel to the decision boundary represent the margins. The blue data points are located above the upper margin line, and the red data points are located below the lower margin line. The text "Support vector machines (SVMs)" is overlaid in the center in a large, bold, red font.

Support vector machines (SVMs)

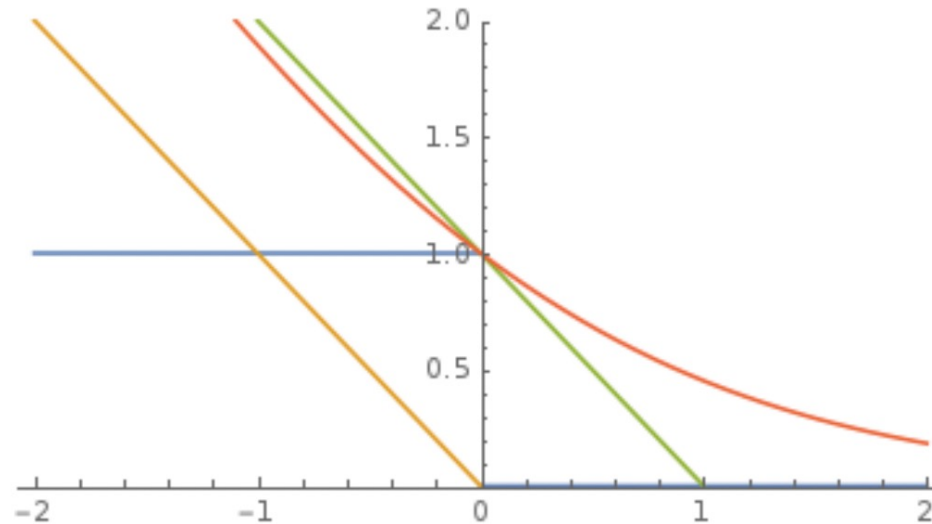
Why study SVM?

- One of the most commonly used classification algorithms
- Allows us to explore the concept of *margins* in classification
- Works well with the kernel trick
- Strong theoretical guarantees

We focus on **binary classification** here.

“Primal” formulation

In one sentence: **linear model with ℓ_2 regularized hinge loss**. Recall:



- **perceptron loss** $\ell_{\text{perceptron}}(z) = \max\{0, -z\} \rightarrow$ Perceptron
- **logistic loss** $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z)) \rightarrow$ logistic regression
- **hinge loss** $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\} \rightarrow$ **SVM**

“Primal” formulation

For a linear model (\mathbf{w}, b) , this means

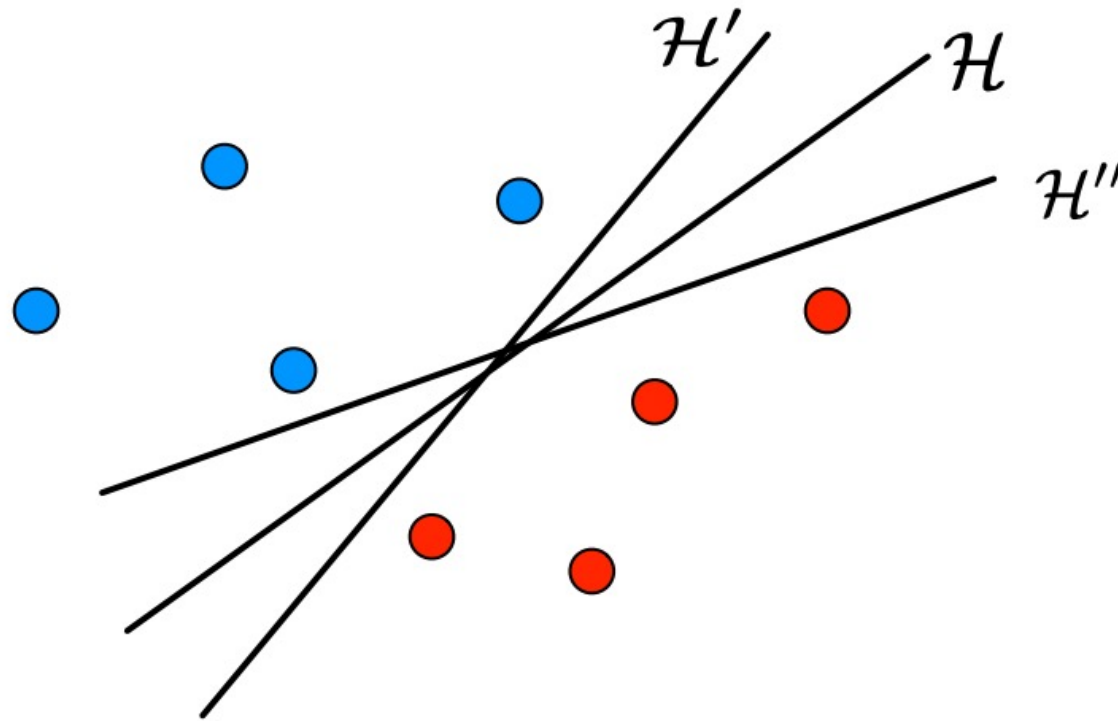
$$\min_{\mathbf{w}, b} \sum_i \max \{0, 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- recall $y_i \in \{-1, +1\}$
- a nonlinear mapping ϕ is applied
- the bias/intercept term b is used explicitly (why?)

So why L2 regularized hinge loss?

Geometric motivation: separable case

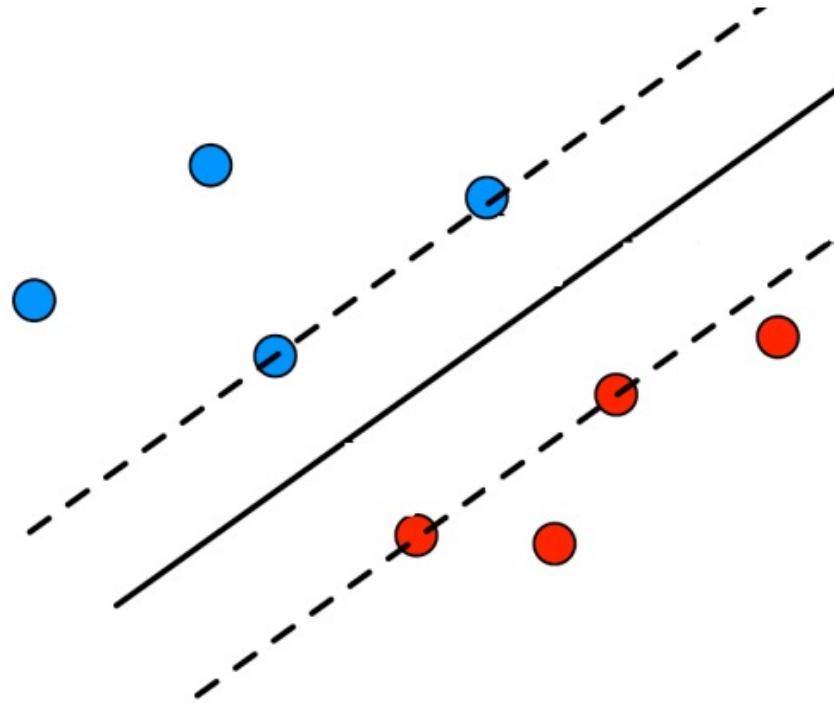
When data is **linearly separable**, there are infinitely many hyperplanes with zero training error:



Which one should we choose?

Geometric motivation: separable case

The further away the separating hyperplane is from the datapoints, the better.



Margins: formalizing this intuition

Distance to hyperplane

What is the **distance** from a point \mathbf{x} to a hyperplane $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$?

Assume the **projection** is $\mathbf{x} - \beta \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, then

$$0 = \mathbf{w}^T \left(\mathbf{x} - \beta \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right) + b = \mathbf{w}^T \mathbf{x} - \beta \|\mathbf{w}\| + b$$

and thus $\beta = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|_2}$.

Therefore the distance is

$$\frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}$$

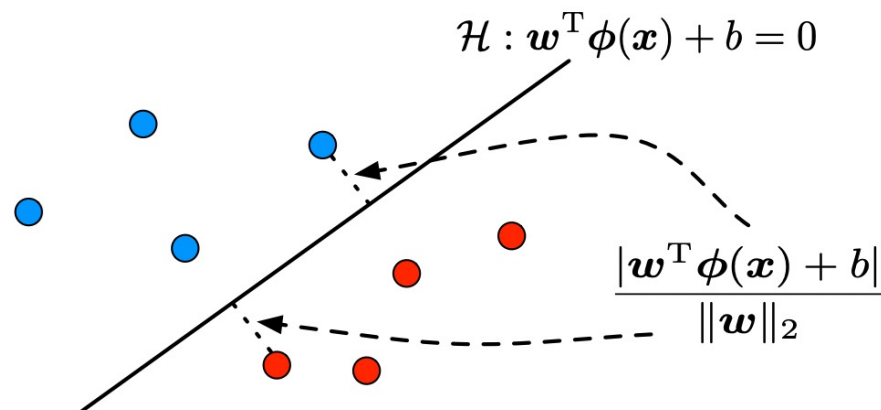
For a hyperplane that correctly classifies (\mathbf{x}, y) , the distance becomes

$$\frac{y(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|_2}$$

Maximizing margin

Margin: the *smallest* distance from all training points to the hyperplane

$$\text{MARGIN OF } (\mathbf{w}, b) = \min_i \frac{y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)}{\|\mathbf{w}\|_2}$$



The intuition “**the further away the better**” translates to solving

$$\max_{\mathbf{w}, b} \min_i \frac{y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)}{\|\mathbf{w}\|_2} = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_i y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)$$

Rescaling

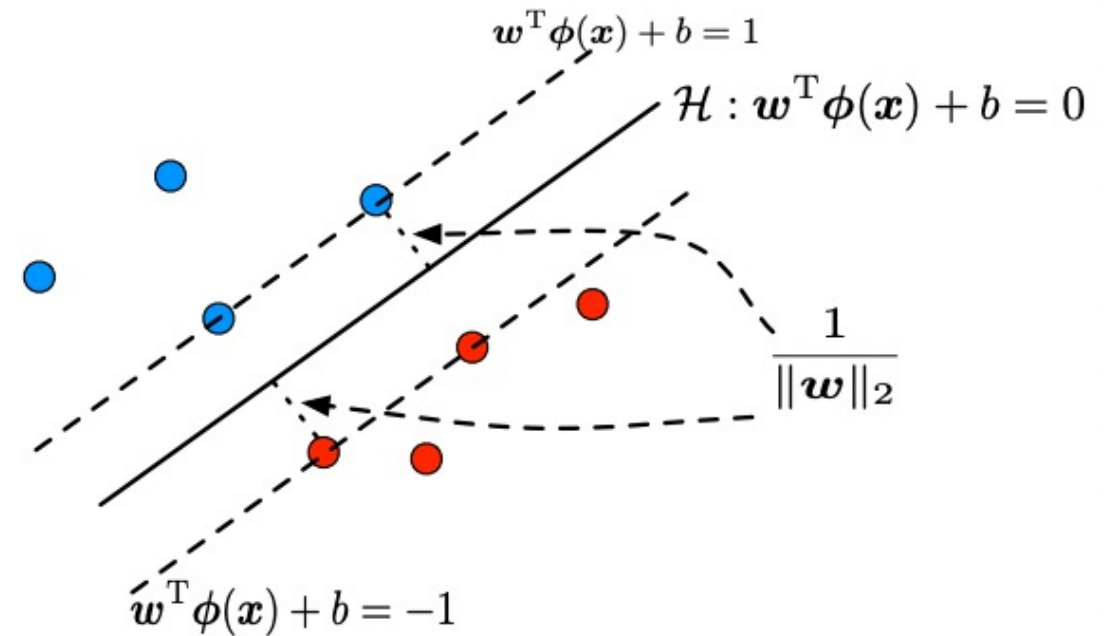
Note: rescaling (\mathbf{w}, b) does not change the hyperplane.

We can thus always scale (\mathbf{w}, b) s.t. $\min_i y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) = 1$

The margin then becomes

MARGIN OF (\mathbf{w}, b)

$$\begin{aligned} &= \frac{1}{\|\mathbf{w}\|_2} \min_i y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \\ &= \frac{1}{\|\mathbf{w}\|_2} \end{aligned}$$



Summary for separable data

For a separable training set, we aim to solve

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{s.t.} \quad \min_i y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) = 1$$

This is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad \forall i \in [n] \end{aligned}$$

SVM is thus also called *max-margin* classifier. The constraints above are called *hard-margin* constraints.

General non-separable case

If data is not linearly separable, the previous constraint

$$y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad \forall i \in [n]$$

is obviously *not feasible*.

To deal with this issue, we relax them to **soft-margin** constraints:

$$y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n]$$

where we introduce **slack variables** $\xi_i \geq 0$.

SVM Primal Formulation

We want ξ_i to be as small as possible too. The objective becomes

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n] \end{aligned}$$

where C is a hyperparameter to balance the two goals.

Equivalent form

The formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \leq \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n] \end{aligned}$$

is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} = \xi_i, \quad \forall i \in [n] \end{aligned}$$

Equivalent form

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} = \xi_i, \quad \forall i \in [n] \end{aligned}$$

is equivalent to

$$\min_{\mathbf{w}, b} C \sum_i \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} + \frac{1}{2} \|\mathbf{w}\|_2^2$$

and

$$\min_{\mathbf{w}, b} \sum_i \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

with $\lambda = 1/C$. *This is exactly minimizing ℓ_2 regularized hinge loss!*

Optimization

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \leq \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n] \end{aligned}$$

- it is a convex (in fact, a **quadratic**) problem
- thus can apply any convex optimization algorithms, e.g. SGD
- there are **more specialized and efficient** algorithms
- but usually we apply kernel trick, which requires solving the *dual problem*