# CSCI 567: Machine Learning

Vatsal Sharan
Spring 2026

Lecture 5, February 13

# Administrivia

- HW2 is out, due in less than 2 weeks.
- Exam 1 in 3 weeks, more details next week

Recap

# Ensuring generalization

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

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

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

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.
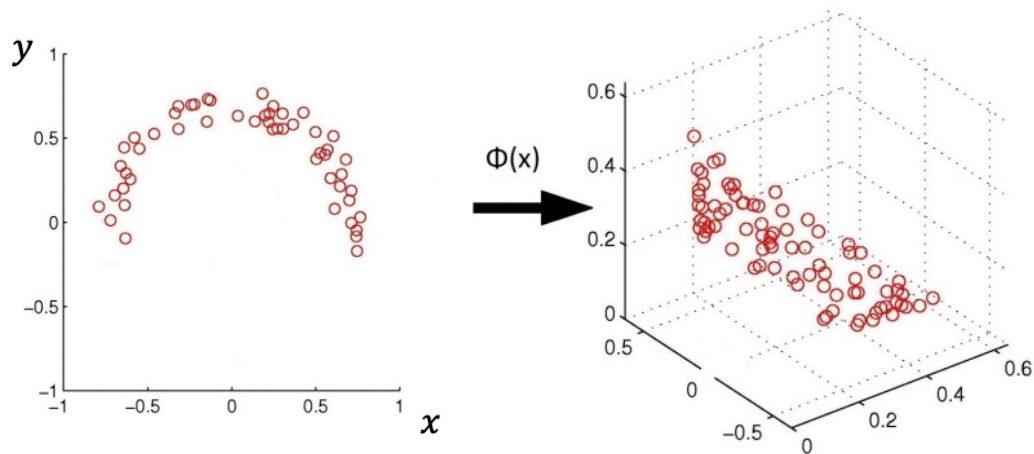
# Beyond linear models: nonlinearly transformed features

**1. Use a nonlinear mapping**

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

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** $d = 1$

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

Learning a linear model in the new space
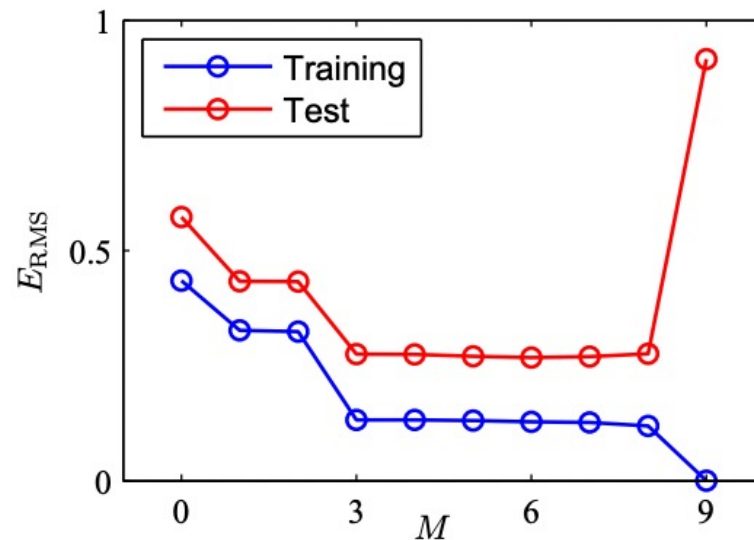= learning an *M-degree polynomial model* in the original space

# 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?

See Colab notebook
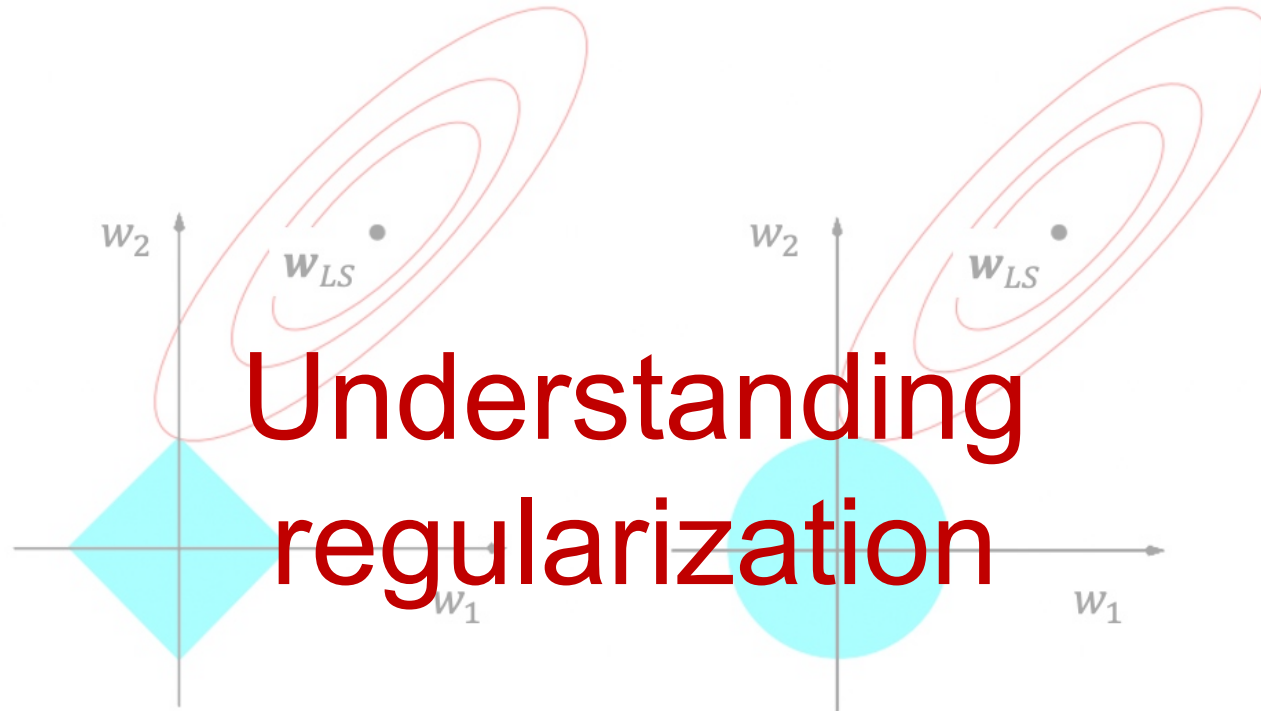
# Preventing overfitting: **Regularization**

**Regularized linear regression**: new objective

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

Goal: find $\boldsymbol{w}^* = \text{argmin}_w\, G(\boldsymbol{w})$

- $\psi : \mathbb{R}^d \to \mathbb{R}^+$ is the *regularizer*

  - measure how complex the model $\boldsymbol{w}$ is, penalize complex models
  - common choices: $\|\boldsymbol{w}\|_2^2$, $\|\boldsymbol{w}\|_1$, etc.

- $\lambda > 0$ is the *regularization coefficient*

  - $\lambda = 0$, no regularization
  - $\lambda \to +\infty$, $\boldsymbol{w} \to \text{argmin}_w\, \psi(\boldsymbol{w})$
  - i.e. control **trade-off** between training error and complexity

Understanding regularization

# $\ell_2$ **regularization**: penalizing large weights

$\ell_2$ **regularization,** $\psi(\boldsymbol{w}) = \|\boldsymbol{w}\|_2^2$:

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

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

$$\Rightarrow \left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X} + \lambda\boldsymbol{I}\right)\boldsymbol{w} = \boldsymbol{X}^{\mathrm{T}}\boldsymbol{y}$$

$$\Rightarrow \boldsymbol{w}^* = \left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X} + \lambda\boldsymbol{I}\right)^{-1}\boldsymbol{X}^{\mathrm{T}}\boldsymbol{y}$$

Linear regression with $\ell_2$ regularization is also known as **ridge regression**.

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

# Encouraging sparsity: $\ell_0$ regularization

Continuing from the frequentist view, having small norm is one possible structure to impose on the model. Another very common one is **sparsity**.

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

E.g. $w = [\,1\,,\,0,-1,0,0.2,0,0\,]$ is 3-sparse

# Encouraging sparsity: $\ell_0$ regularization

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

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.

# Encouraging sparsity: $\ell_0$ regularization

**Sparsity of $w$:** Number of non-zero coefficients in $w$. Same as $||\mathbf{w}||_0$
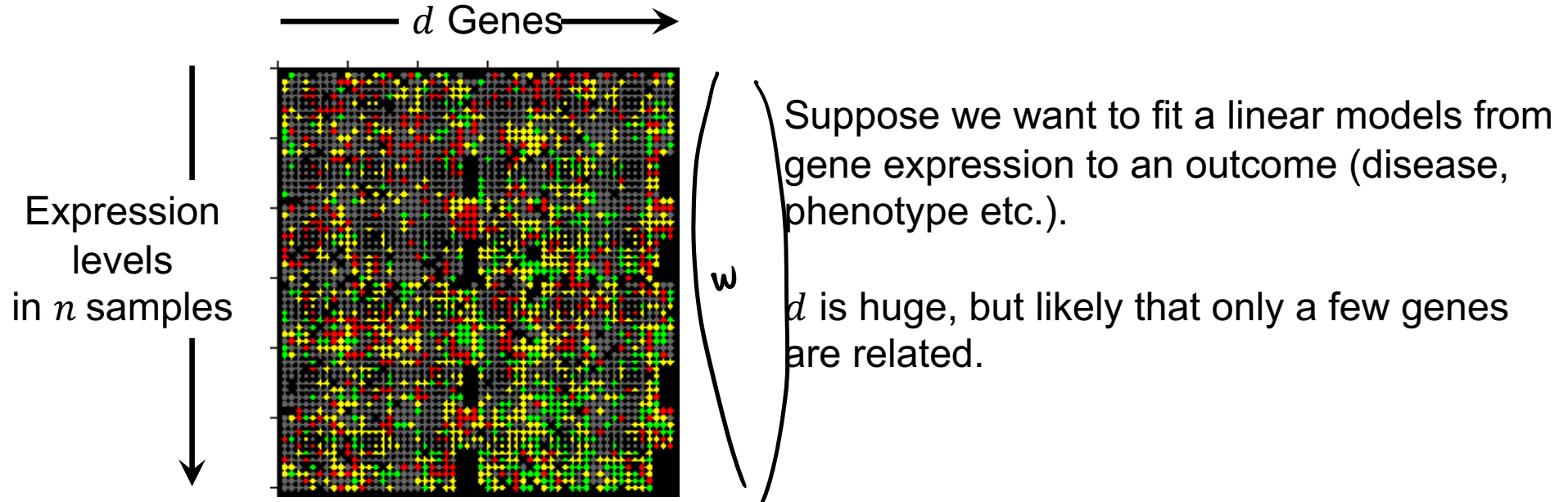
Advantage:

- 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.



$\overrightarrow{\phantom{--}d \text{ Genes}\phantom{--}}$

Expression levels in $n$ samples

Suppose we want to fit a linear models from gene expression to an outcome (disease, phenotype etc.).

$d$ is huge, but likely that only a few genes are related.

# Encouraging sparsity: $\ell_0$ regularization

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

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.

E.g.     $w = [\mathbf{1.5}, 0, -\mathbf{1.1}, 0, \mathbf{0.25}, 0, 0\,]$ is more interpretable than,
         $w = [1, 0.2, -1.3, 0.15, 0.2, 0.05, 0.12\,]$

For a sparse model, it could be easier to understand the model. It is also easier to verify whether the features which have a high weight have a relation with the outcome (they are not spurious artifacts of the data).

# Encouraging sparsity: $\ell_0$ regularization

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

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.

# $\ell_0$ regularization: The good and the bad

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

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

Good: Need less data to learn

Suppose weights in $w$ are in $\{-1, 0, 1\}$

How many such $s$-sparse vectors are there in $d$ dimensions?

Answer: $\binom{d}{s} \cdot 2^s$

# $\ell_0$ regularization: The good and the bad

How much data to learn?

About $\cong \log(|F|)$ samples to learn

$$= \log\left(\binom{d}{s} \cdot 2^s\right) = \log\left(\binom{d}{s}\right) + \log 2^s$$
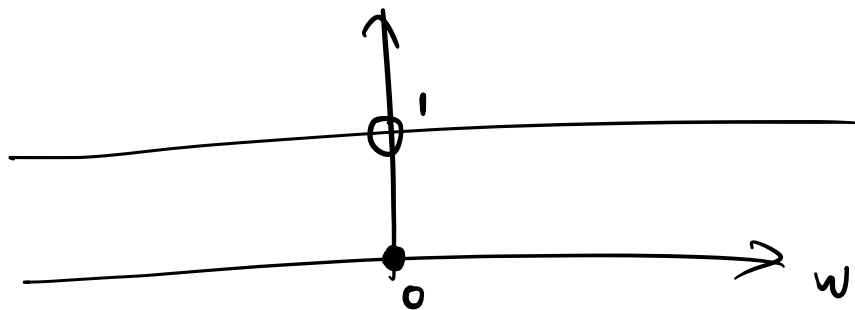
$$\leq \log d^s + \log 2^s$$

$$= s \log d + s \log 2$$

$$\left(\binom{d}{s} \leq d^s\right)$$

In contrast, without s-sparsity constraint need about $d$ samples in $d$-dimensions.

$\therefore$ If $s \ll d$, need much less data to generalize!

# $\ell_0$ regularization: The good and the bad

**Bad**: $\|w\|_0$ is highly discontinuous



GD has no hope!

# $\ell_1$ regularization as a proxy for $\ell_0$ regularization

Choose $\psi(\boldsymbol{w}) = \|\boldsymbol{w}\|_1. \ = \sum_{\partial=1}^{d} |w_{\partial}|$

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

There is theory which says that under some appropriate conditions, doing $\ell_1$ regularization has the same effect as if we did $\ell_0$ regularization, i.e. we get sparsity, and have the same data requirement as if we did $\ell_0$ regularization!

$\|w\|_1$ is convex $\Rightarrow$ GD/ SGD to solve $G(w)$ !

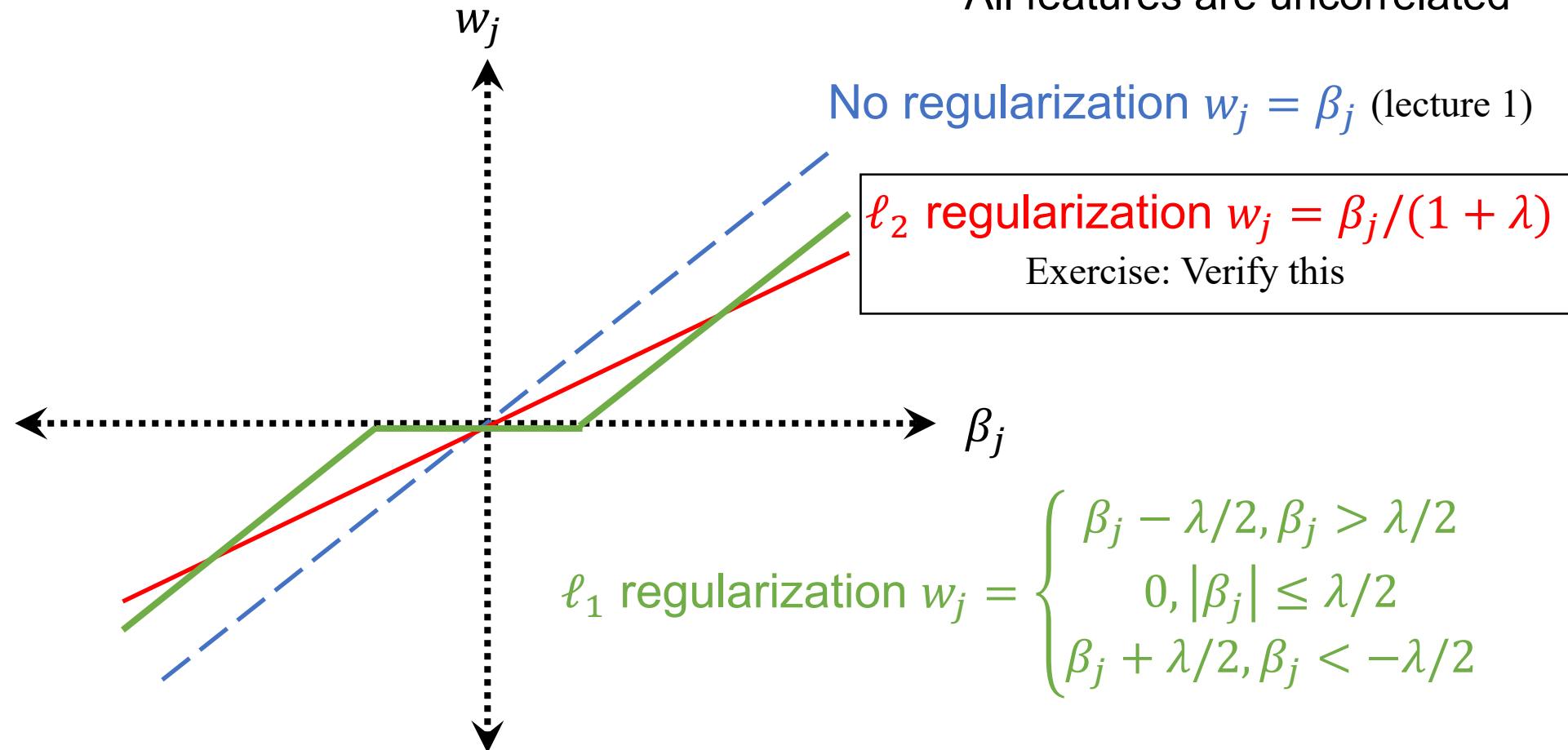# Diving deeper: $\ell_1$ and $\ell_2$ regularization for the "isotropic" case

Isotropic case ($X^T X = I$)
$X^T_{(j)}$ is $j$-th row of $X^T$ ($j$-th feature of all datapoints)
Let $\beta_j = X^T_{(j)} y$ , (correlation of $j$-th feature with label)

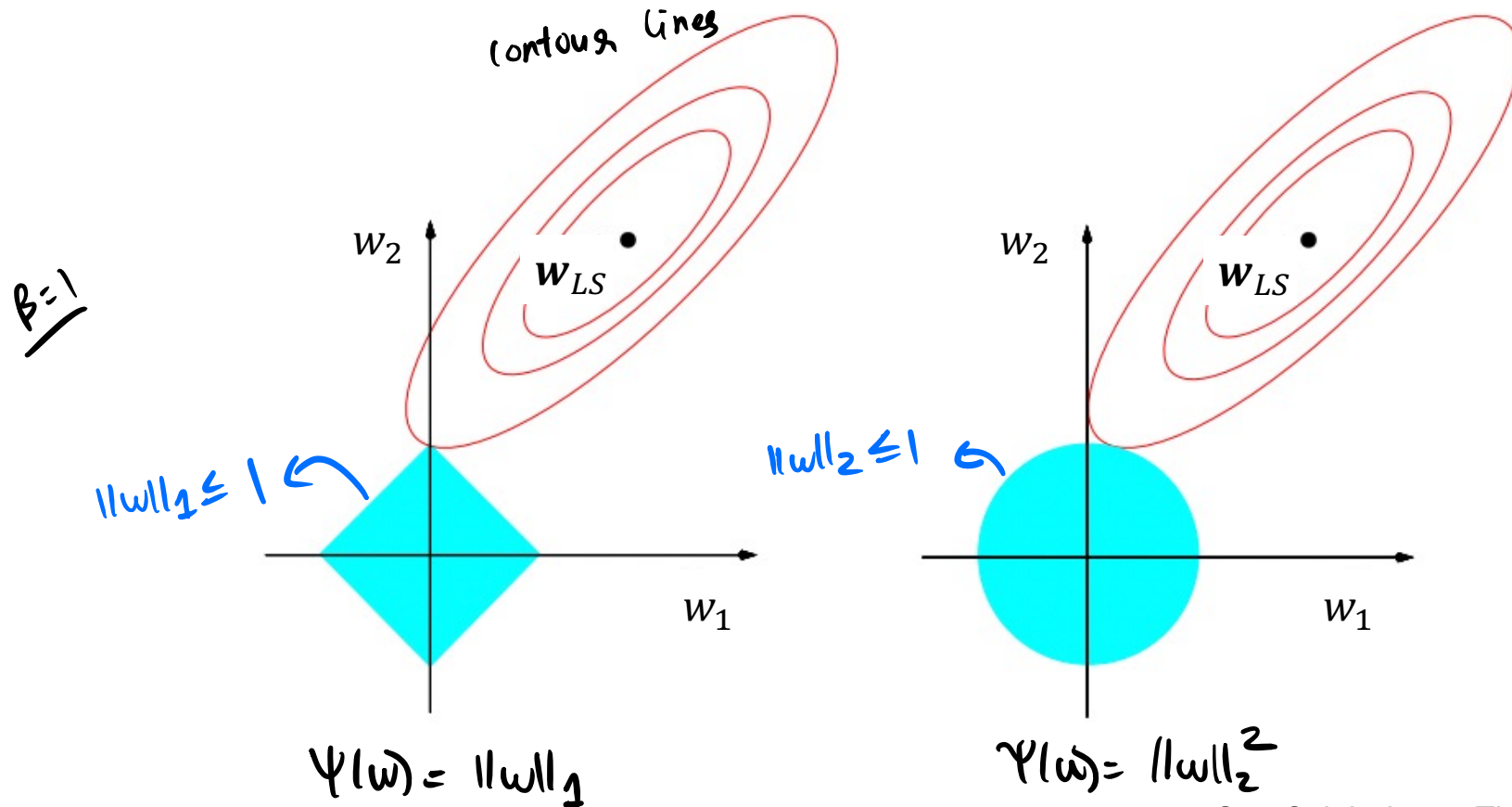Isotropic case ($X^T X = I$)
- Take all features to be mean 0
- All features have variance 1
- All features are uncorrelated



No regularization $w_j = \beta_j$ (lecture 1)

$\ell_2$ regularization $w_j = \beta_j/(1 + \lambda)$

Exercise: Verify this

$\ell_1$ regularization $w_j = \begin{cases} \beta_j - \lambda/2, \beta_j > \lambda/2 \\ 0, |\beta_j| \leq \lambda/2 \\ \beta_j + \lambda/2, \beta_j < -\lambda/2 \end{cases}$

# Why does $\ell_1$ regularization encourage sparse solutions?

Optimization problem: $\text{argmin}_w \text{ RSS}(w)$ , subject to $\psi(w) \leq \beta$



$\beta = 1$

contour lines

$w_2$

$w_{LS}$

$\|w\|_1 \leq 1$

$\psi(w) = \|w\|_1$

$w_2$

$w_{LS}$

$\|w\|_2 \leq 1$

$\psi(w) = \|w\|_2^2$

$w_1$

$w_1$

# Implicit regularization

So far, we explicitly added a $\psi(\boldsymbol{w})$ term to our objective function to regularize.
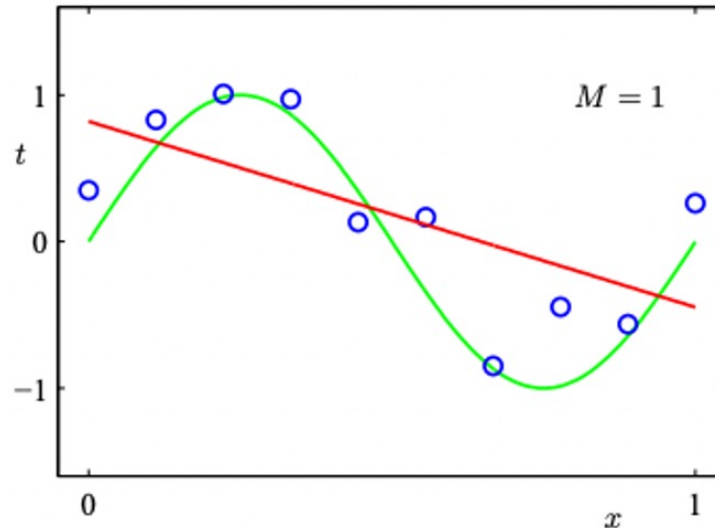
In many cases, the optimization algorithm we use can themselves act as regularizers, favoring some solutions over others.

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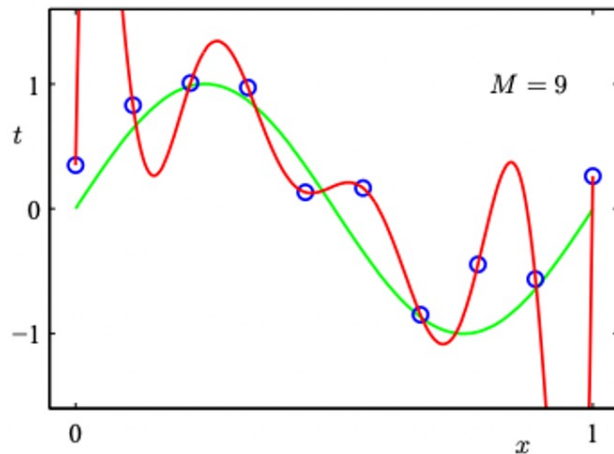
A model whose complexity is too *small* for the task will *underfit*. This is a model with a large bias because the model's accuracy will not improve even if we add a lot of training data.



sin(x) fitting example we saw in Lec 3

# 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).



sin(x) fitting example we saw in Lec 3

Input Space

$\phi$

Feature Space

**Kernels**

# Motivation

Recall the nonlinear function map for linear regression:

1. **Use a nonlinear mapping**

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

   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 $\phi : \mathbb{R}^d \to \mathbb{R}^M$ (for linear regression, and much more broadly).

# Regularized least squares

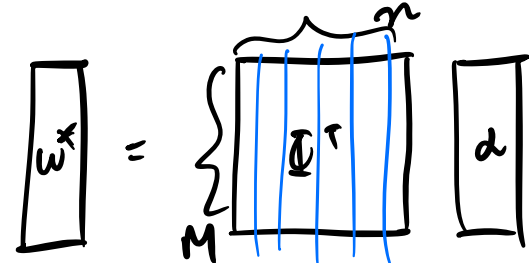Let's continue with regularized least squares with non-linear basis:

$$\boldsymbol{w}^* = \operatorname*{argmin}_{\boldsymbol{w}} F(\boldsymbol{w})$$

$$= \operatorname*{argmin}_{\boldsymbol{w}} \left( \|\boldsymbol{\Phi}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{w}\|_2^2 \right)$$

$$= \left( \boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{\Phi} + \lambda \boldsymbol{I} \right)^{-1} \boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{y}$$

$$\boldsymbol{\Phi} \underset{\in \mathbb{R}^{n+M}}{=} \begin{pmatrix} \phi(\boldsymbol{x}_1)^{\mathrm{T}} \in \mathbb{R}^M \\ \phi(\boldsymbol{x}_2)^{\mathrm{T}} \\ \vdots \\ \phi(\boldsymbol{x}_n)^{\mathrm{T}} \end{pmatrix}, \quad \boldsymbol{y} \underset{\in \mathbb{R}^n}{=} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

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 $F(\boldsymbol{w}) = \|\boldsymbol{\Phi}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \lambda\|\boldsymbol{w}\|_2^2$ to be $\boldsymbol{0}$:

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

we know

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

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 $\boldsymbol{\alpha}$, the prediction of $\boldsymbol{w}^*$ on a new example $\boldsymbol{x}$ is

$$\boldsymbol{w}^{*\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i \boldsymbol{\phi}(\boldsymbol{x}_i)^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}) \quad = \left( \sum_{i=1}^{n} \alpha_i \, \phi(x_i) \right)^{\mathrm{T}} \phi(x)$$

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

Kernel methods are exactly about computing inner products *without explicitly computing $\boldsymbol{\phi}$.*

But we need to figure out what $\boldsymbol{\alpha}$ is first!

# Solving for $\alpha$, Step 1: Kernel matrix

Plugging in $\boldsymbol{w} = \boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\alpha}$ into $F(\boldsymbol{w})$ gives

$$= \left(\boldsymbol{\Phi}^{\mathrm{T}}\alpha\right)^{\mathrm{T}}\left(\boldsymbol{\Phi}^{\mathrm{T}}\alpha\right) = \alpha^{\mathrm{T}}\boldsymbol{\Phi}\,\boldsymbol{\Phi}^{\mathrm{T}}\alpha$$

$$H(\boldsymbol{\alpha}) = F(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\alpha})$$

$$= \|\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\alpha} - \boldsymbol{y}\|_2^2 + \lambda\|\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\alpha}\|_2^2$$

$$= \|\boldsymbol{K}\boldsymbol{\alpha} - \boldsymbol{y}\|_2^2 + \lambda\boldsymbol{\alpha}^{\mathrm{T}}\boldsymbol{K}\boldsymbol{\alpha} \qquad (\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} \in \mathbb{R}^{n \times n})$$

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

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

# Kernel matrix: Example

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

**Gram/Kernel matrix**

$$\boldsymbol{K} = \begin{pmatrix} \phi(x_1)^{\mathrm{T}}\phi(x_1) & \phi(x_1)^{\mathrm{T}}\phi(x_2) & \phi(x_1)^{\mathrm{T}}\phi(x_3) \\ \phi(x_2)^{\mathrm{T}}\phi(x_1) & \phi(x_2)^{\mathrm{T}}\phi(x_2) & \phi(x_2)^{\mathrm{T}}\phi(x_3) \\ \phi(x_3)^{\mathrm{T}}\phi(x_1) & \phi(x_3)^{\mathrm{T}}\phi(x_2) & \phi(x_3)^{\mathrm{T}}\phi(x_3) \end{pmatrix}$$

$$= \begin{pmatrix} 4 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 4 \end{pmatrix}$$

# Kernel matrix vs Covariance matrix

| | dimensions | entry $(i,j)$ | property |
|---|---|---|---|
| $\mathbf{\Phi\Phi}^{\mathrm{T}}$ | $n \times n$ | $\phi(x_i)^{\mathsf{T}} \phi(x_j)$ | Both are symmetric |
| $\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}$ | $M \times M$ | $\displaystyle\sum_{k=1}^{n} \underbrace{\phi(x_k)_i} \cdot \phi(x_k)_j$ | & positive semi-definite (psd) |

ith co-ordinate of k-th datapoint

why are they psd?

Any matrix $A = VV^{\mathsf{T}}$ is psd

$\left( x^{\mathsf{T}} A x = x^{\mathsf{T}} VV^{\mathsf{T}} x = (V^{\mathsf{T}}x)^{\mathsf{T}} (V^{\mathsf{T}}x) = \| V^{\mathsf{T}}x \|_2^2 \geqslant 0 \right)$

# Solving for $\alpha$, Step 2: Minimize the dual

Minimize (the so-called *dual formulation*)

$$H(\boldsymbol{\alpha}) = \|\boldsymbol{K\alpha} - \boldsymbol{y}\|_2^2 + \lambda \boldsymbol{\alpha}^{\mathrm{T}} \boldsymbol{K\alpha}$$

Setting the derivative to $\mathbf{0}$ we have

$$\mathbf{0} = (\boldsymbol{K}^2 + \lambda \boldsymbol{K})\boldsymbol{\alpha} - \boldsymbol{Ky} = \boldsymbol{K}\left(\underbrace{(\boldsymbol{K} + \lambda \boldsymbol{I})\boldsymbol{\alpha} - \boldsymbol{y}}_{\text{set to } 0}\right)$$

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

$$\boldsymbol{w}^* = \boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\alpha} = \boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{K} + \lambda \boldsymbol{I})^{-1}\boldsymbol{y}$$

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

# Comparing two solutions

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

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

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:**

$$(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi} + \lambda I)^{-1}\mathbf{\Phi}^\mathsf{T}y$$
$$= (\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi} + \lambda I)^{-1}\mathbf{\Phi}^\mathsf{T}(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T} + \lambda I)(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T} + \lambda I)^{-1}y$$
$$= (\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi} + \lambda I)^{-1}(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T} + \lambda\mathbf{\Phi}^\mathsf{T})(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T} + \lambda I)^{-1}y$$
$$= (\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi} + \lambda I)^{-1}(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi} + \lambda I)\mathbf{\Phi}^\mathsf{T}(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T} + \lambda I)^{-1}y$$
$$= \mathbf{\Phi}^\mathsf{T}(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T} + \lambda I)^{-1}y$$

# The kernel trick

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

First, computing $(\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}} + \lambda\mathbf{I})^{-1}$ can be more efficient than computing $(\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi} + \lambda\mathbf{I})^{-1}$ when $n \leq M$.　$\in \mathbb{R}^{n+n} \to O(n^3)$ time　$\in \mathbb{R}^{M+M}$
$\to O(M^3)$

More importantly, computing $\boldsymbol{\alpha} = (\boldsymbol{K} + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$ 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 $\phi(\boldsymbol{x})^T\phi(\boldsymbol{x}')$.*

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

# The kernel trick: Example 1

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

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

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

$$\phi(x)^{\mathrm{T}}\phi(x') = {x_1}^2{x_1'}^2 + 2x_1 x_2 x_1' x_2' + {x_2}^2{x_2'}^2$$
$$= (x_1 x_1' + x_2 x_2')^2 = (x^{\mathrm{T}} x')^2$$

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

$\phi : \mathbb{R}^d \to \mathbb{R}^{2d}$ is parameterized by $\theta$:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \qquad \phi_\theta(x) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_d) \\ \sin(\theta x_d) \end{pmatrix}$$

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

$$\phi_\theta(x)^{\mathrm{T}} \phi_\theta(x') = \sum_{m=1}^{d} \cos(\theta x_m) \cos(\theta x'_m) + \sin(\theta x_m) \sin(\theta x'_m)$$

$$= \sum_{m=1}^{d} \cos(\theta(x_m - x'_m)) \qquad \text{(trigonometric identity)}$$

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 $\phi_\theta$, define $\phi_L : \mathbb{R}^d \to \mathbb{R}^{2d(L+1)}$ for some integer $L$:

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

$\theta$ varies from

$\left( 0, \frac{2\pi}{L}, \frac{4\pi}{L}, \cdots, L\cdot\frac{2\pi}{L} \right)$

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

$$\phi_L(x)^{\mathrm{T}}\phi_L(x') = \sum_{\ell=0}^{L} \phi_{\frac{2\pi\ell}{L}}(x)^{\mathrm{T}}\phi_{\frac{2\pi\ell}{L}}(x')$$

$$= \sum_{\ell=0}^{L}\sum_{m=1}^{d} \cos\left( \frac{2\pi\ell}{L}(x_m - x'_m) \right)$$

# The kernel trick: Example 4

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

$$\phi_\infty(\boldsymbol{x})^{\mathrm{T}}\phi_\infty(\boldsymbol{x}') = \int_0^{2\pi} \sum_{m=1}^{d} \cos(\theta(x_m - x'_m))\, d\theta$$

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

$$\frac{2\pi \ell}{L} = \theta$$

$$\int_0^{2\pi} \cos(x\theta)\, d\theta$$

$$= \frac{\sin(x\theta)}{x} \Big|_0^{2\pi}$$

Again, a simple function of the original features.

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

# Kernel functions

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

$$k(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}')$$

Examples we have seen

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^{\mathrm{T}} \boldsymbol{x}')^2$$

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

# Using kernel functions

Choosing a nonlinear basis $\phi$ 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:

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

$\nearrow \ \phi(x_i)^{\mathrm{T}} \phi(x_n)$

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

- useful for proving that a function is not a kernel

# Examples which are not kernels

Function

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

is *not a kernel*, why?

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

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

must be positive semidefinite, *but is it?*

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \text{is} \quad \text{not} \quad \text{psd. why ?}$$

$$\begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} -1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = -2 !$$

# Properties of kernels

For any function $f : \mathbb{R}^d \to \mathbb{R}$, $k(\boldsymbol{x}, \boldsymbol{x}') = f(\boldsymbol{x})f(\boldsymbol{x}')$ 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**: $k_1(\cdot, \cdot)k_2(\cdot, \cdot)$ $\to$ HW 2

- **exponential**: $e^{k(\cdot, \cdot)}$

- $\cdots$

Verify using the definition of kernel!

feature mapping $\phi(x)$

s.t. $k(x, x') = \phi(x)^\top \phi(x')$

$\phi(x) = f(x)$

what is $\phi$ ?

$\phi_1$: map for $k_1$

$\phi_2$: map for $k_2$

$\phi'$: map for $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$

Exercise : find $\phi'$

# Popular kernels

**Polynomial kernel**

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

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

*polynomial in original i·p.*

What is the corresponding $\phi$?

$$c = 0, \quad M = 2 \quad , \quad k(x, x') = \left(x^{\mathrm{T}}x'\right)^2$$

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

# Popular kernels

**Gaussian kernel or Radial basis function (RBF) kernel**

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

$$\|x\|_2^2 + \|x'\|_2^2 - 2\,x^T x'$$

What is the corresponding $\phi$?

is this a kernel ?

$$k(x, x') = \exp\left(-\frac{\|x\|_2^2}{2\sigma^2}\right) \exp\left(-\frac{\|x'\|_2^2}{2\sigma^2}\right) \boxed{\exp\left(\frac{x^T x'}{\sigma^2}\right)}$$

$$k(x, x') = f(x)\, f(x')$$

$$\text{where} \quad f(x) = \exp\left(-\frac{\|x\|_2^2}{2\sigma^2}\right)$$

transformation   for   product

# Popular kernels

**Gaussian kernel or Radial basis function (RBF) kernel**

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

What is the corresponding $\phi$?

$$e^{x} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

$$\exp\left(\frac{x^T x'}{\sigma^2}\right) = 1 + \frac{x^T x'}{\sigma^2} + \frac{1}{2!}\frac{(x^T x')^2}{(\sigma^2)^2} + \frac{1}{3!}\frac{(x^T x')^3}{(\sigma^2)^3} + \dots$$

each of these is
polynomial kernel

# Popular kernels

Appropriate kernels have also been developed for tasks like Natural Language Processing where inputs are discrete.

For two strings $s_1$ and $s_2$ and some parameter $t$,

$$k_t(s_1, s_2) = \text{Number of sub-strings of length } t \text{ which appear in both } s_1 \text{ and } s_2.$$

For e.g. if $t = 1$,

$$k_t(\text{`machine'}, \text{`learning'}) = 4.$$

What is the corresponding $\phi$?

# Prediction with kernels

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

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

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 $\boldsymbol{w}^*$ 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 $\text{sign}(\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}))$. Data may become linearly separable in the feature space!

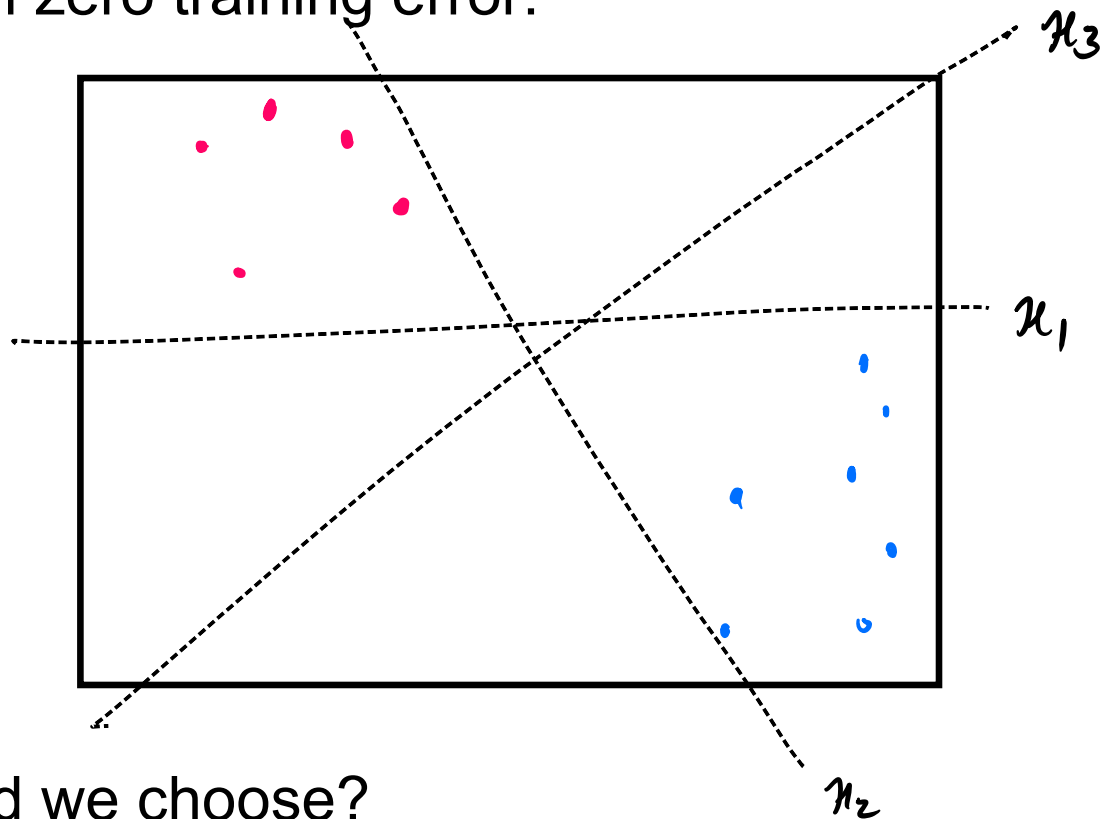# Support vector machines (SVMs)

# 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.

The *function class for SVMs is a linear function on a feature map $\phi$ applied to the datapoints*: $\text{sign}(\boldsymbol{w}^{\text{T}}\phi(\boldsymbol{x}) + b)$. Note, the bias term $b$ is taken separately for SVMs, you'll see why.

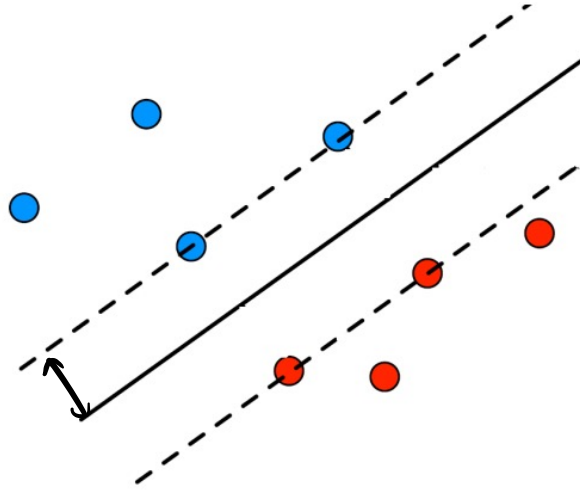# Margins: separable case, geometric intuition

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



Which one should we choose?

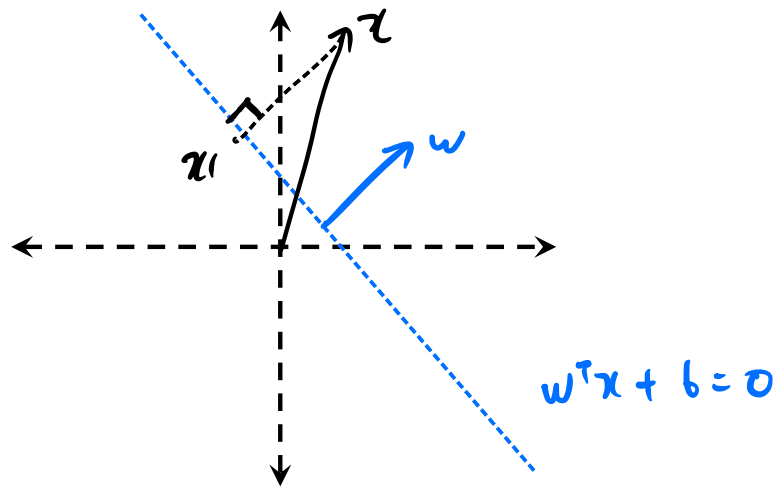# Margins: separable case, geometric intuition

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



Margin for linearly separable data: Distance from the hyperplane to the point closest to hyperplane

# Formalizing geometric intuition: Distance to hyperplane

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



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

$$\left( \text{for any } \alpha, \ \alpha \, \text{sign}(\alpha) = |\alpha| \right)$$

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

Therefore the distance is $\|x - x'\|_2 = |\beta| = \frac{|w^T x + b|}{\|w\|_2}$.

$$\text{sign}(w^T x + b) = y$$

For a hyperplane that correctly classifies $(x, y)$, the distance becomes $\frac{y(w^T x + b)}{\|w\|_2}$.