

CS-GY 6923: Lecture 9

Kernel Methods, Support Vector Machines

NYU Tandon School of Engineering, Prof. Christopher Musco

NON-LINEAR METHODS

- Many previous methods studied (regression, logistic regression) are considered linear methods. They make predictions based on $\langle \mathbf{x}, \beta \rangle$ – i.e. based on weighted sums of features.
- In the next part of the course we move on to non-linear methods. Specifically, kernel methods and neural networks.
- Both are very closely related to feature transformations, which was a technique we already saw for using linear methods to learn non-linear concepts.

RECALL: k -NEAREST NEIGHBOR METHOD

k -NN algorithm: a simple but powerful baseline for classification.

Training data: $(\underline{x}_1, \underline{y}_1), \dots, (\underline{x}_n, \underline{y}_n)$ where $\underline{y}_1, \dots, \underline{y}_n \in \{1, \dots, q\}$.

Classification algorithm:

Given new input \underline{x}_{new} ,

SOME similarity function (sim)

- Compute $\underline{sim}(\underline{x}_{new}, \underline{x}_1), \dots, \underline{sim}(\underline{x}_{new}, \underline{x}_n)$.¹
- Let $\underline{x}_{j_1}, \dots, \underline{x}_{j_k}$ be the training data vectors with highest similarity to \underline{x}_{new} .
- Predict \underline{y}_{new} as $majority(\underline{y}_{j_1}, \dots, \underline{y}_{j_k})$.

¹ $sim(\underline{x}_{new}, \underline{x}_i)$ is any chosen similarity function, like $1 - \|\underline{x}_{new} - \underline{x}_i\|_2$.

k -NEAREST NEIGHBOR METHOD

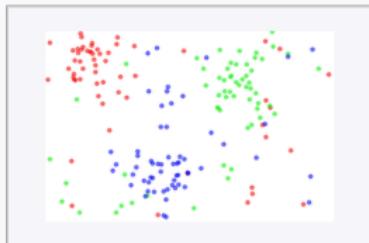


Fig. 1. The dataset.



Fig. 2. The 1NN classification map.

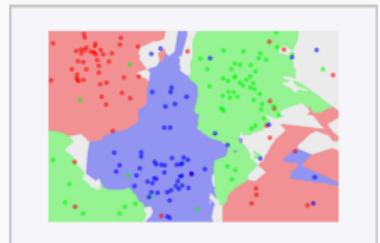


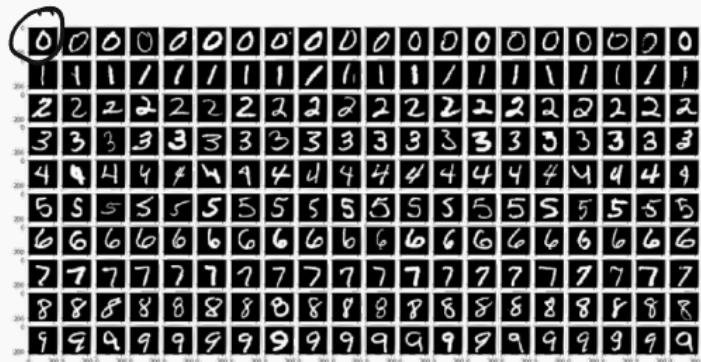
Fig. 3. The 5NN classification map.

- { Smaller k , more complex classification function.
- { Larger k , more robust to noisy labels.

Works remarkably well for many datasets.

MNIST IMAGE DATA

Especially good for large datasets with lots of repetition. Works well on MNIST for example. 95% Accuracy out-of-the-box.

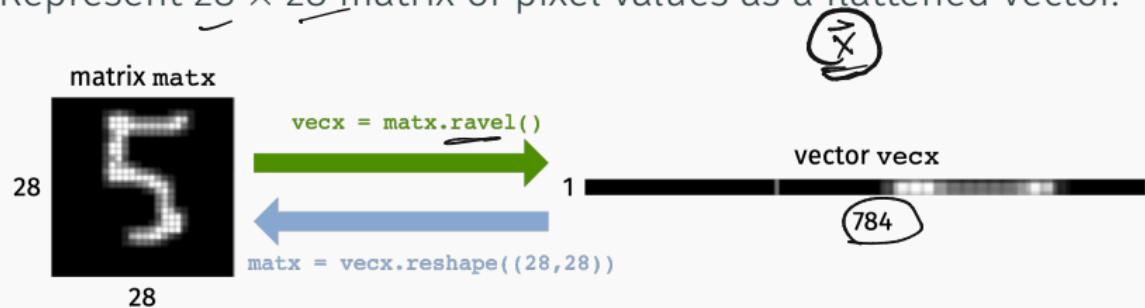


Let's look into this example a bit more...

MNIST IMAGE DATA

Each pixel is number from [0, 1]. 0 is black, 1 is white.

Represent 28×28 matrix of pixel values as a flattened vector.



```
xmat = np.array([[1,2,3],[4,5,6],[7,8,9]])
```

```
array([[1, 2, 3],  
       [4, 5, 6],  
       [7, 8, 9]])
```

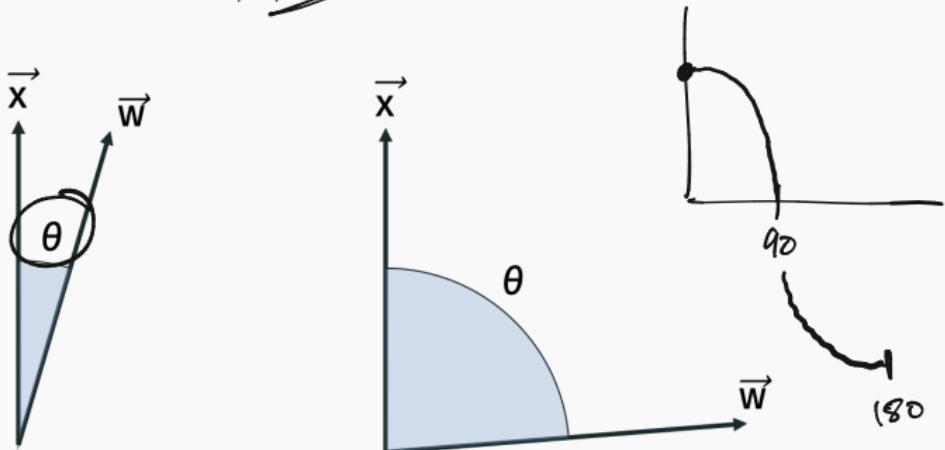
```
xvec = xmat.ravel()
```

```
array([1, 2, 3, 4, 5, 6, 7, 8, 9])
```

INNER PRODUCT SIMILARITY

Given data vectors $\mathbf{x}, \mathbf{w} \in \mathbb{R}^d$, the inner product $\langle \mathbf{x}, \mathbf{w} \rangle$ is a natural similarity measure.

$$\langle \mathbf{x}, \mathbf{w} \rangle = \sum_{i=1}^d x_i w_i = \underline{\cos(\theta)} \underline{\|\mathbf{x}\|_2} \underline{\|\mathbf{w}\|_2}.$$



Also called “cosine similarity”.

INNER PRODUCT SIMILARITY

$$(x - \omega)^\top (x - \omega) = \underline{x^\top x} + \underline{\omega^\top \omega} - \underline{2x^\top \omega}$$

Connection to Euclidean (ℓ_2) Distance:

$$\|\underline{x - w}\|_2^2 = (\|\underline{x}\|_2^2) + (\|\underline{w}\|_2^2) - \underline{2\langle x, w \rangle}$$

If all data vectors has the same norm, the pair of vectors with largest inner product is the pair with smallest Euclidean distance.

INNER PRODUCT FOR MNIST

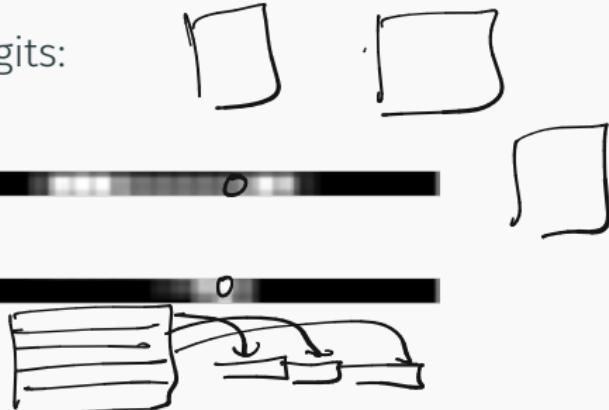
Inner product between MNIST digits:

x
black

\vec{x}
white

\vec{w}
white

$$\langle \underline{x}, w \rangle = \sum_{i=1}^{28} \sum_{j=1}^{28} \text{matx}[i,j] \cdot \text{matw}[i,j].$$

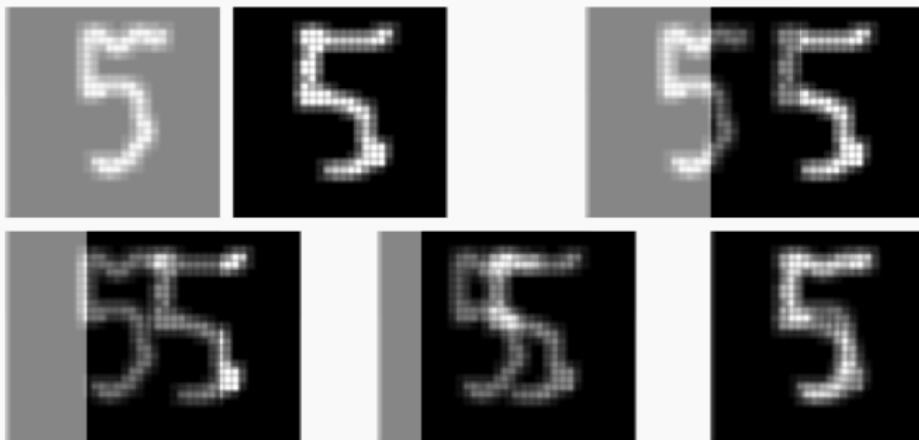


Inner product similarity is higher when the images have large pixel values (close to 1) in the same locations. I.e. when they have a lot of overlapping white/light gray pixels.

$B^{(0)}$

INNER PRODUCT FOR MNIST

Visualizing the inner product between two images:



Images with high inner product have a lot of overlap.

K-NN ALGORITHM ON MNIST



Most similar images during k -nn search, $k = 9$:

A 5x10 grid of handwritten digits. A large curly brace on the left side groups the first five rows. Handwritten annotations are present: a circled '8' in the fourth row, a circled '5' in the fourth row, and a circled '3' in the third row. The digits are as follows:

2	2	2	2	2	2	2	2	2	2
6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	3	7	7
8	8	8	8	5	8	8	8	5	8
0	0	6	0	6	6	6	0	6	6

ANOTHER VIEW ON LINEAR CLASSIFICATION

$0, \dots, q$

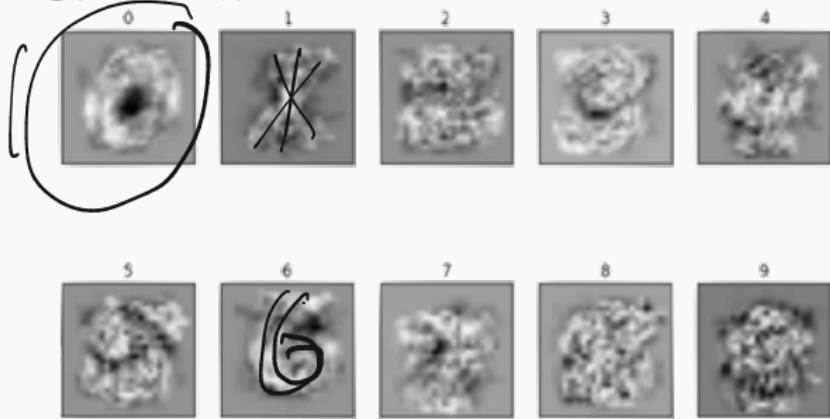
One-vs.-all or Multiclass Cross-entropy Classification with Logistic Regression:

- Learn q classifiers with parameters $\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(q)}$.
- Given x_{new} compute $\langle x_{new}, \beta^{(1)} \rangle, \dots, \langle x_{new}, \beta^{(q)} \rangle$
- Predict class $y_{new} = \arg \max_i \langle x_{new}, \beta^{(i)} \rangle$.

If each x is a vector with $28 \times 28 = 784$ entries than each $\beta^{(i)}$ also has 784 entries. Each parameter vector can be viewed as a 28×28 image.

MATCHED FILTER

Visualizing $\beta^{(0)}, \dots, \beta^{(q)}$:

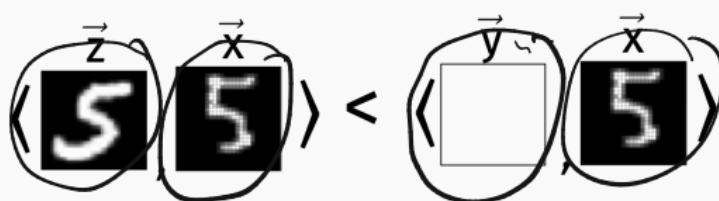


Logistic regression classification rule: For an input , compute inner product similarity with all weight matrices and choose most similar one.

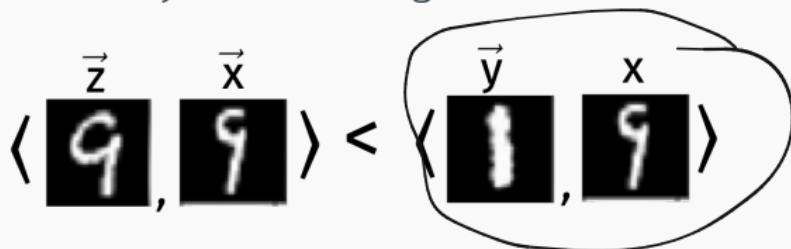
In contrast to k -NN, only needs to compute similarity with q items instead of n . Much faster classification.

DIVING INTO SIMILARITY

Often the inner product **does not make sense** as a similarity measure between data vectors. Here's an example (recall that smaller inner product means less similar):



But clearly the first image is more similar.



Here's a more realistic scenario.

KERNEL FUNCTIONS: ALTERNATIVE MEASURES OF SIMILARITY

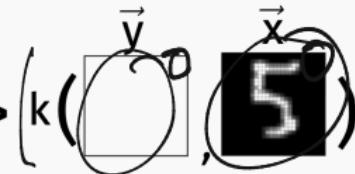
A kernel function $k(x, y)$ is simply a similarity measure between data points.

$$k(x, y) = \begin{cases} \text{large if } x \text{ and } y \text{ are similar.} \\ \cancel{\text{close to 0 if } x \text{ and } y \text{ are different.}} \\ \text{smaller} \end{cases}$$

Example: The Radial Basis Function (RBF) kernel, aka the Gaussian kernel:

$$k(x, y) = e^{-\|\underline{x}-\underline{y}\|_2^2/\sigma^2}$$

for some scaling factor σ .

$$k(\vec{z}, \vec{x}) > k(\vec{y}, \vec{x})$$


KERNEL FUNCTIONS: A NEW MEASURE OF SIMILARITY

Lots of kernel functions functions involve transformations of $\langle \mathbf{x}, \mathbf{y} \rangle$ or $\|\mathbf{x} - \mathbf{y}\|_2$:

- Gaussian RBF Kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_2^2/\sigma^2}$
- Laplace Kernel: $k(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x}-\mathbf{y}\|_2/\sigma}$
- Polynomial Kernel: $k(\mathbf{x}, \mathbf{y}) = (\underbrace{\langle \mathbf{x}, \mathbf{y} \rangle + 1}_q)^q.$ }

But you can imagine much more complex similarity metrics.

We will see one on the next problem set tailored to digit/letter recognition.

HOW TO USE A KERNEL FUNCTION?

For k -nearest neighbors, can easily replace inner product with whatever similarity function you want.

For logistic regression, it is less clear how to do so.

HOW TO USE A KERNEL FUNCTION?

Logistic Regression Loss:

$$L(\beta^{(1)}, \dots, \beta^{(q)}) = - \underbrace{\left(\sum_{i=1}^n \left(\sum_{\ell=1}^q \mathbb{1}[y_i = \ell] \cdot \log \frac{e^{\langle \beta^{(\ell)}, x_i \rangle}}{\sum_{j=1}^q e^{\langle \beta^{(j)}, x_i \rangle}} \right) \right)}_{f(\beta^\ell, x_j)}$$

Loss inherently involves inner product between each $\beta^{(j)}$ and each data vector x_i .

Solution: Only work with similarity metrics that can be expressed as inner products.

KERNEL FUNCTIONS FROM FEATURE TRANSFORMATION

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

$$1 \dots n$$

2

\vdots

n

G

$k(x_i, x_j) = G_{ij}$

Kernel Gram Matrix

A positive semidefinite (PSD) kernel is any similarity function with the following form:

$$k(x, w) = \langle \phi(x)^\top \phi(w) \rangle = \langle \phi(x), \phi(w) \rangle$$

where $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is a some feature transformation function.

$$\phi([]) \rightarrow []$$

~~all zero~~

KERNEL FUNCTIONS AND FEATURE TRANSFORMATION

Example: Degree 2 polynomial kernel, $k(\mathbf{x}, \mathbf{w}) = (\mathbf{x}^T \mathbf{w} + 1)^2 = \phi(\mathbf{x})^T \phi(\mathbf{w})$

$$\mathbf{x}^T \begin{bmatrix} A \\ \end{bmatrix} \mathbf{x} \geq 0$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$(x^T w + 1)^2 = (x_1 w_1 + x_2 w_2 + x_3 w_3 + 1)^2$$

$$= 1 + 2x_1 w_1 + 2x_2 w_2 + 2x_3 w_3 + x_1^2 w_1^2 + x_2^2 w_2^2 + x_3^2 w_3^2 + 2x_1 w_1 x_2 w_2 + 2x_1 w_1 x_3 w_3 + 2x_2 w_2 x_3 w_3$$

$$= \phi(\mathbf{x})^T \phi(\mathbf{w}).$$

$$\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_3 \\ x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1 x_2 \\ \sqrt{2}x_1 x_3 \\ \sqrt{2}x_2 x_3 \\ 0 \end{bmatrix}$$

$$\phi(\mathbf{w}) = \begin{bmatrix} 1 \\ \frac{f_1 w_1}{f_2 w_1} \\ \vdots \\ \frac{f_2 w_1}{f_2 w_3} \\ \vdots \\ \sqrt{2} w_1 w_3 \\ f_2 w_1 w_3 \\ 0 \end{bmatrix}$$

KERNEL FUNCTIONS AND FEATURE TRANSFORMATION

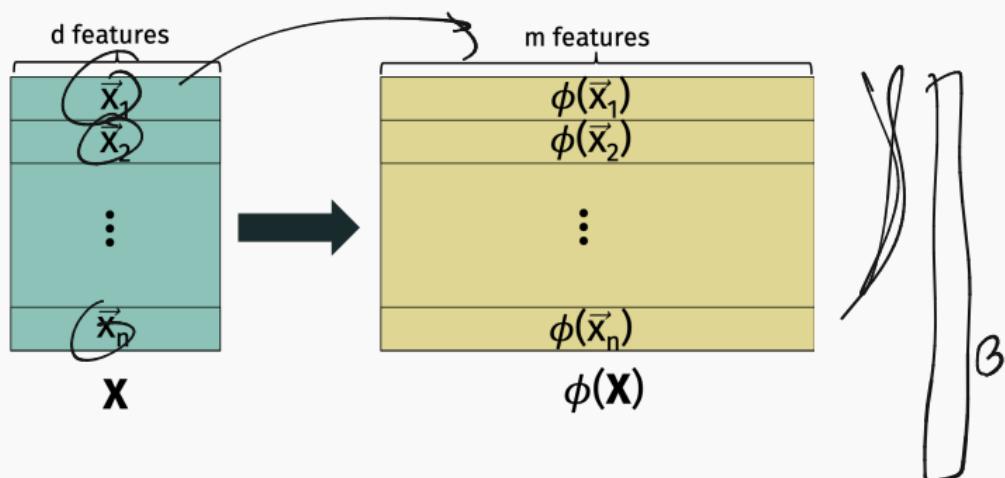
Not all similarity metrics are positive semidefinite (PSD), but all of the ones we saw earlier are:

- { Gaussian RBF Kernel: $k(x, y) = e^{-\|x-y\|_2^2/\sigma^2}$ } $O(d)$
 - { Laplace Kernel: $k(x, y) = e^{-\|x-y\|_2/\sigma}$ } $O(d)$
 - { Polynomial Kernel: $k(x, y) = \underbrace{\langle x, y \rangle + 1}_A$ } $O(d)$
- And there are many more... \ddot{d}^b

KERNEL FUNCTIONS AND FEATURE TRANSFORMATION

Feature transformations \iff new similarity metrics.

To work with the similarity $k(\cdot, \cdot)$ in place of the inner product $\langle \cdot, \cdot \rangle$, it suffices to replace every data point $\mathbf{x}_1, \dots, \mathbf{x}_n$ by $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)$.



There are two major issues with this:

- While $\phi(\mathbf{x})$ is sometimes simple and explicit. **More often, it is not.** We might be able to show a kernel is PSD without easily being able to write down $\phi(\mathbf{x})$.
- Transform dimension m is often very large: e.g. $m = O(d^q)$ for a degree q polynomial kernel. For many kernels (e.g. the Gaussian kernel) m is actually *infinite*.

So doing the feature transformation explicitly would have very high computational cost.

REPARAMETERIZATION TRICK

Any vector β representing x^α can be reparameterized as x^α for class α .

For simplicity, let's just consider the binary cross entropy/logistic regression loss:

$$x^\alpha \leftarrow h(x, \beta)^\top$$

$$\text{min}_\beta - \sum_{j=1}^n y_j \log(h(X\beta)_j) + (1 - y_j) \log(1 - h(X\beta)_j)$$

s.t. β represents x^α

$$\text{where } h(z) = \frac{1}{1+e^{-z}}$$



REPARAMETERIZATION TRICK

Reminder from linear algebra: Without loss of generality, can assume that β lies in the row span of X . all vectors $\in \mathbb{R}^d$ that can be written as $\sum_{i=1}^n c_i x_i$ for coefficients c_1, \dots, c_n

So for any $\beta \in \mathbb{R}^d$, there exists a vector $\alpha \in \mathbb{R}^n$ such that:

$$\underbrace{X\beta}_{=Xv} = \underbrace{XX^T\alpha}_{=Xv}$$

$$a_1x_1 + a_2x_2 + \dots + a_nx_n$$

$$\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \quad \beta \quad = \quad \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \quad \begin{matrix} x_1 & x_2 & \dots & x_n \end{matrix} \quad v \quad \begin{matrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{matrix}$$

$$X(\beta) = X(v + e) = Xv + Xe$$

$$= Xv + 0 \\ = Xv$$

$S = \text{row space of } S$, write β as $v + e$

REPARAMETERIZATION TRICK

$$h(X\beta) \quad h(XX^T\alpha) \quad \underline{X\beta} \rightarrow XX^T\alpha$$

Logistic Regression Equivalent Formulation: Given data matrix $X \in \mathbb{R}^{n \times d}$ and binary label vector $y \in \{0, 1\}^n$ for class i , find $\underline{\alpha} \in \mathbb{R}^n$ to minimize the loss:

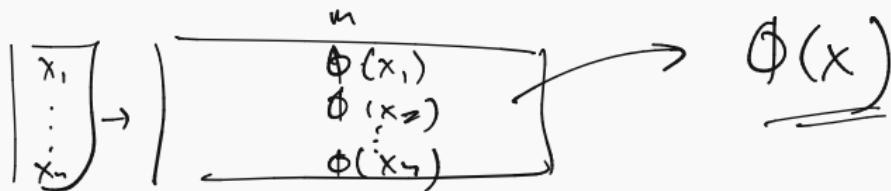
$$\min_{\alpha} - \sum_{j=1}^n y_j \log(h(XX^T\alpha)_j) + (1 - y_j) \log(1 - h(XX^T\alpha)_j)$$

Can still be minimized via gradient descent:

$$\underbrace{\nabla L(\alpha)}_{\nabla L(\alpha) = X X^T (h(X X^T \alpha) - y)} = X X^T (h(X X^T \alpha) - y).$$

$$X^T(h(X\beta) - y) \quad \nabla L(\alpha) = X X^T (h(X X^T \alpha) - y)$$

REPARAMETERIZATION TRICK



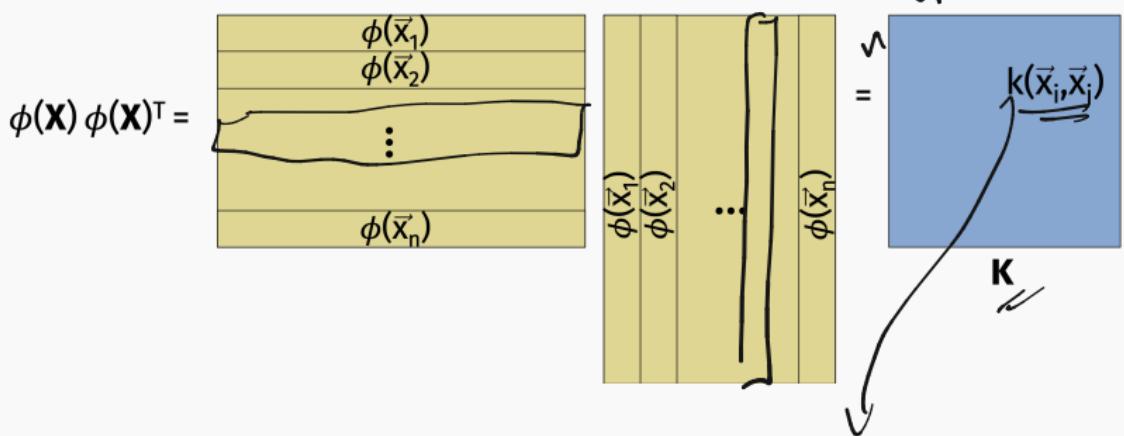
If we use a non-linear data transformation ϕ (corresponding to a PSD kernel), then the loss is:

$$-\sum_{j=1}^n y_j \log(h(\phi(X)\phi(X)^T\alpha)_j) + (1 - y_j) \log(1 - h(\phi(X)\phi(X)^T\alpha)_j)$$

The diagram illustrates the components of the loss function. The term $y_j \log(h(\phi(X)\phi(X)^T\alpha)_j)$ is associated with a matrix X of size $n \times m$, and the term $(1 - y_j) \log(1 - h(\phi(X)\phi(X)^T\alpha)_j)$ is associated with a matrix X of size $n \times n$. Both terms involve the function h and the parameter vector α .

KERNEL MATRIX

$K = \phi(X)\phi(X)^T$ is called the (kernel Gram matrix)

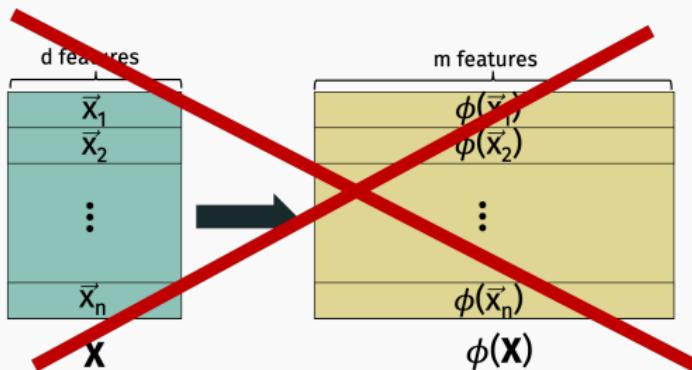


$$\langle \phi(x_i), \phi(x_j) \rangle = k(x_i, x_j)$$

KERNEL TRICK

We never need to actually compute $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)$ explicitly!

- For training we just need the kernel matrix \mathbf{K} which requires computing $k(\mathbf{x}_i, \mathbf{x}_j)$ for all i, j .



We can always work with a finite sized $n \times n$ matrix.

KERNEL TRICK

Take away:

$$\boxed{\text{ }} \boxed{\beta} \rightarrow \sum_{i=1}^n \phi(x_i) \alpha_i$$

- Logistic regression can be combined with any positive semidefinite kernel matrix, and the model can be trained in time independent of the transform dimension m .

(Prediction can also be done efficiently. For a new input x_{new} , we need to compute:

$$\langle \phi(x_{new}), \beta \rangle = \langle \phi(x_{new}), \phi(X)\alpha \rangle = \sum_{j=1}^n \alpha_j \langle \phi(x_{new}), \phi(x_j) \rangle.$$

\downarrow

$= \langle \phi(x_{new}) | N \rangle$

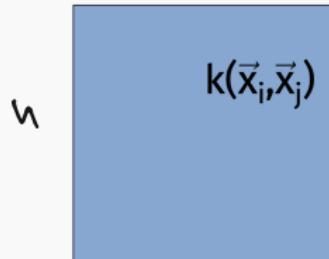
Each term in the sum $\langle \phi(x_{new}), \phi(x_j) \rangle = k(x_{new}, x_j)$ can be computed without explicit feature transformation.²

²Note that this does require computing the kernel inner product with potentially all examples in the training data. More on this shortly.

$$\text{Classification rule } \mathbb{1}[(\langle x, \beta \rangle) > 0] = \mathbb{1}[\text{sum} > 0]$$

BEYOND THE KERNEL TRICK

The kernel matrix \mathbf{K} is still $n \times n$ though which is huge when the size of the training set n is large. Has made the kernel trick less appealing in some modern ML applications.



A blue square representing a matrix \mathbf{K} . In the top-left corner of the square, there is a small white italicized letter 'n'. In the center of the square, there is a white mathematical expression $k(\vec{x}_i, \vec{x}_j)$.

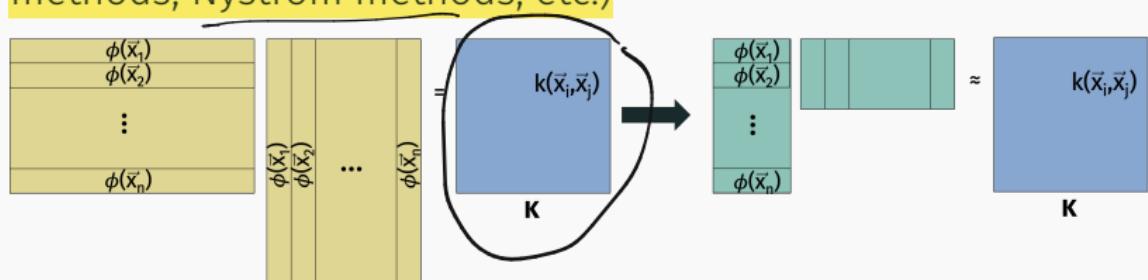
\mathbf{K}

There is an inherent quadratic dependence on n in the computational and space complexity of kernel methods.

- 10,000 data points → runtime scales as $\sim 100,000,000$, \mathbf{K} takes 800MB of space.
- 1,000,000 data points → runtime scales as $\sim 10^{12}$, \mathbf{K} takes 8TB of space.

BEYOND THE KERNEL TRICK

Many algorithmic advances in recent years partially address this computational challenge (random Fourier features methods, Nystrom methods, etc.)



Often based on “reversing” the kernel trick to find a compact feature set that well approximates the kernel.

KERNEL REGRESSION

The kernel trick can also be applied outside of classification.

E.g. to regression:

$$X\beta \rightarrow X\alpha$$

$$\min_{\beta} \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2 \rightarrow \min_{\alpha} \|X\alpha - y\|_2^2 + \lambda \|X\alpha\|_2^2$$

Replace XX^T by kernel matrix K during training.

Prediction:

$$\langle x_{new}, \beta \rangle$$

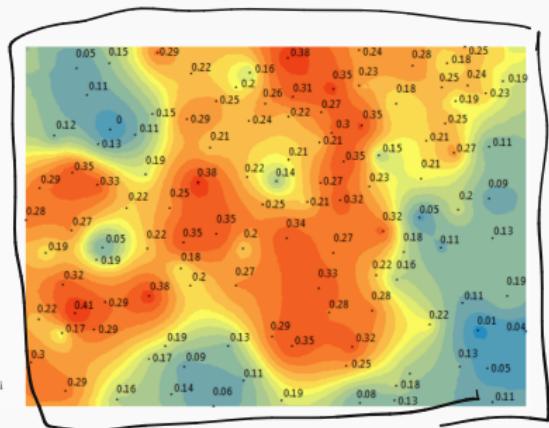
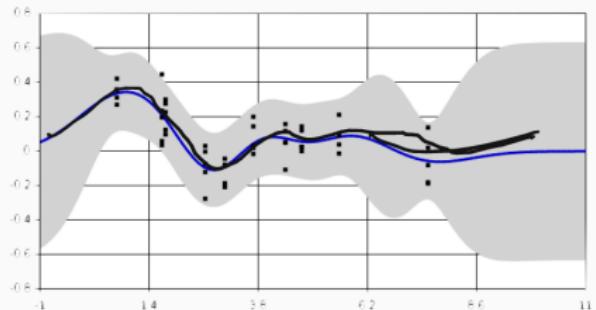
$$y_{new} = \sum_{i=1}^n \alpha_i \cdot k(x_{new}, x_i).$$

Added benefit: Relatively numerically stable. E.g. is a much better option for performing multivariate or even single variate polynomial regression than direct feature expansion.

KERNEL REGRESSION

Kernel regression with non-linear kernels like $e^{-\|x-y\|_2^2}$ is a very important statistical tool, especially when dealing with spatial or temporal data.

Prach Xo
 $\gamma = 10 \text{ pm}$



Also known as Gaussian Process (GP) Regression or Kriging.

Most commonly, Gaussian kernel $k(x, y) = e^{-\|x-y\|_2^2/\sigma^2}$ is used in place of, e.g., polynomial kernel/polynomial regression.

SUPPORT VECTOR MACHINES

Support Vector Machines (SVMs): Another algorithm for finding linear classifiers which is (was?) as popular as logistic regression.

- Can also be combined with kernels.
- Developed from a pretty different perspective.
- But final algorithm is not that different.



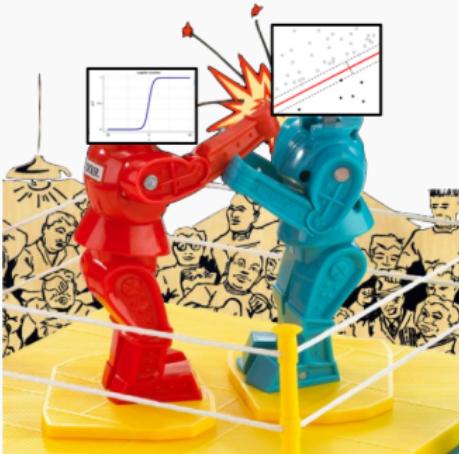
- Invented in 1963 by Alexey Chervonenkis and Vladimir Vapnik. Also founders of VC-theory.
- First combined with non-linear kernels in 1993.

SVM'S VS. LOGISTIC REGRESSION

SVMs are more commonly associated with non-linear kernels. For example, `sklearn`'s SVM classifier (called SVC) has support for non-linear kernels built in by default. Its logistic regression classifier does not.

- Seems to be partially for historical reasons.
- In the early 2000s SVMs were a “hot topic” in machine learning and their popularity persists.
- It is not clear to me if they are better than logistic regression, but honestly the jury is still out...
- There are some computational advantages of using SVMs, and some disadvantages, which is part of the story.

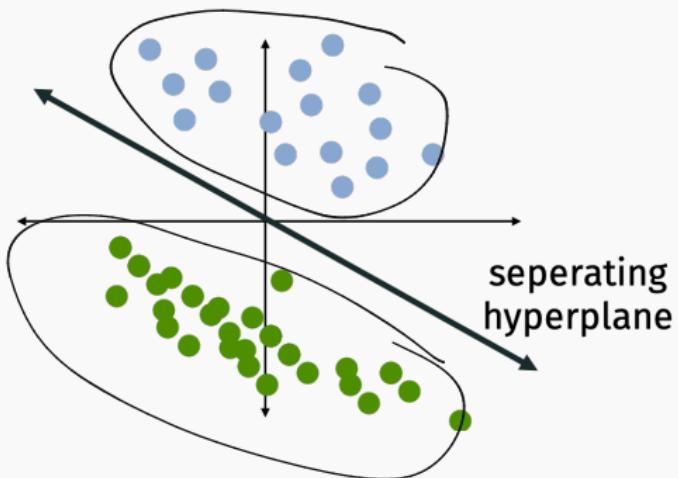
SVM'S VS. LOGISTIC REGRESSION



Next lab: Machina-a-machina comparison of SVMs vs. logistic regression for a MNIST digit classification problem. Which provides better accuracy? Which is faster to train?

LINEARLY SEPARABLE DATA

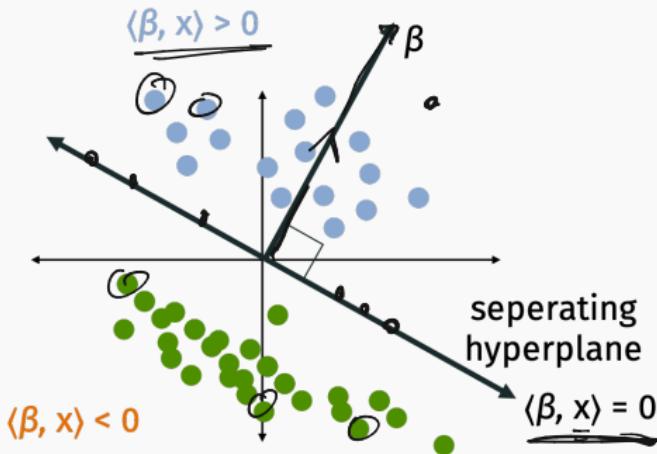
We call a dataset with binary labels linearly separable if it can be perfectly classified with a linear classifier:



This the (realizable) setting we discussed last lecture.

LINEARLY SEPARABLE DATA

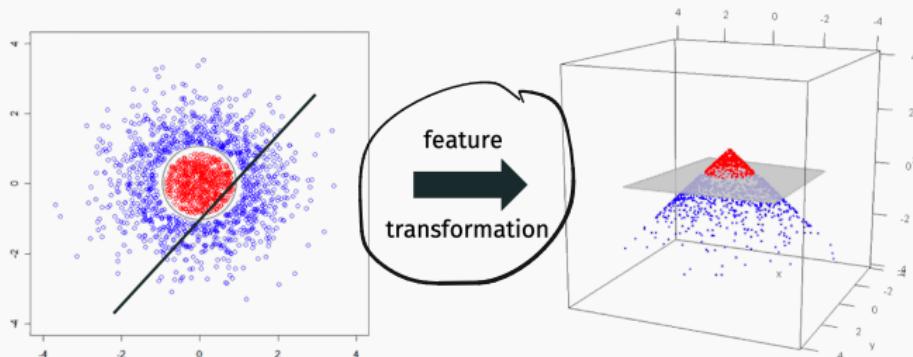
Formally, there exists a parameter β such that $\langle \beta, x \rangle > 0$ for all x in class 1 and $\langle \beta, x \rangle < 0$ for all x in class 0.



Note that if we multiply β by any constant c , $c\beta$ gives the same separating hyperplane because $\langle c\beta, x \rangle = c\langle \beta, x \rangle$.

LINEARLY SEPARABLE DATA

A data set might be linearly separable when using a non-kernel/feature transformation even if it is not separable in the original space.

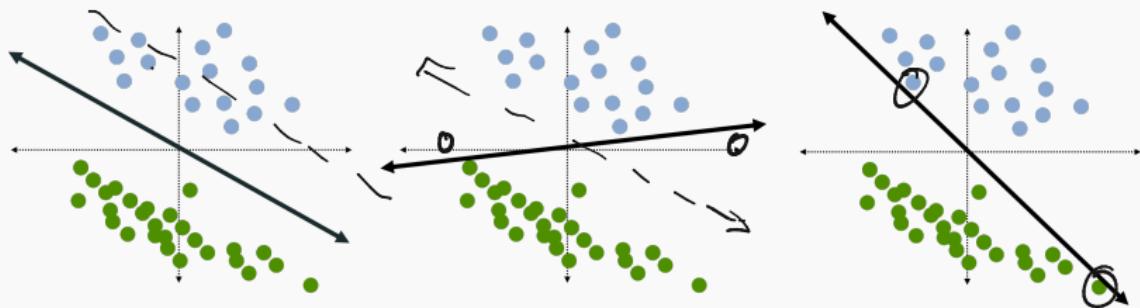


This data is separable when using a degree-2 polynomial kernel. It suffices for $\phi(x)$ to contain x_1^2 and x_2^2 .

When data is linearly separable, we would typically be concerned about over-fitting.

$$\langle \beta, x \rangle > 0$$

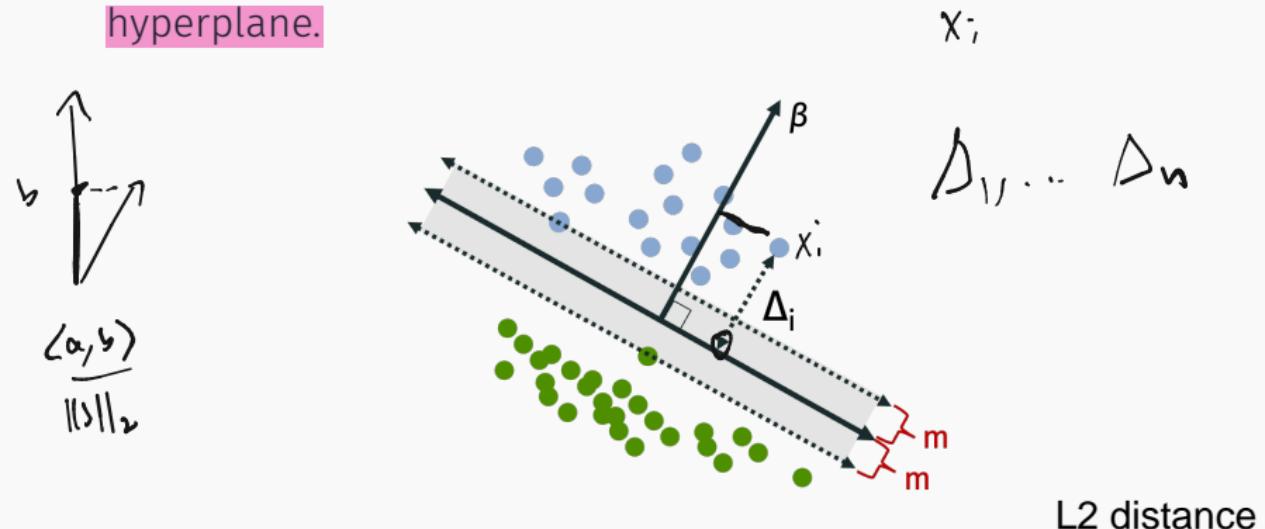
Idea from Vapnik and Chervonenkis: Maybe there is a way to find a classification rule that over-fits on the training data but is still “good” for future data.



There are typically multiple valid separating hyperplanes.
Intuitively, which is best for classifying future data?

MARGIN

The margin m of a separating hyperplane is the minimum ℓ_2 (Euclidean) distance between a point in the dataset and the hyperplane.



$$m = \min_i \Delta_i$$

where

$$\left(\Delta_i = \frac{|\langle x_i, \beta \rangle|}{\|\beta\|_2} \right)$$

We have that $\mathbf{x}_i = \mathbf{v}_i + \mathbf{e}_i$ where \mathbf{v}_i is parallel to $\boldsymbol{\beta}$ and \mathbf{e}_i is perpendicular.

$$\Delta_i = \|\mathbf{v}_i\|_2 = \frac{1}{\|\mathbf{v}_i\|_2} \cdot \langle \mathbf{v}_i, \mathbf{v}_i \rangle = \frac{1}{\|\mathbf{v}_i\|_2} \cdot \frac{\|\mathbf{v}_i\|_2}{\|\boldsymbol{\beta}\|_2} \cdot |\langle \mathbf{v}_i, \boldsymbol{\beta} \rangle| = \frac{|\langle \mathbf{v}_i, \boldsymbol{\beta} \rangle|}{\|\boldsymbol{\beta}\|_2}.$$

Finally, we have that $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle = \langle \mathbf{v}_i, \boldsymbol{\beta} \rangle$ because $\langle \mathbf{e}_i, \boldsymbol{\beta} \rangle = 0$.

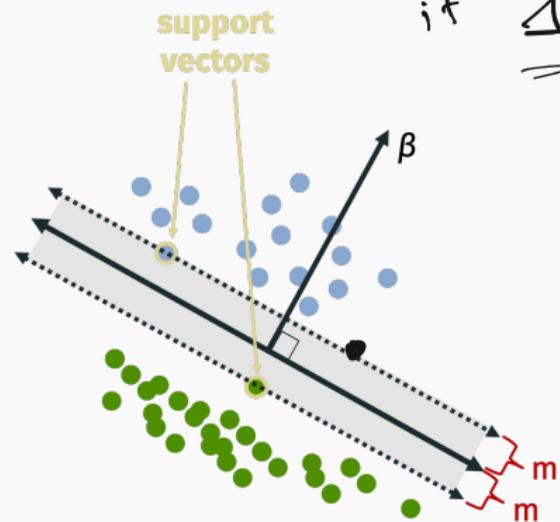
SUPPORT VECTOR

A support vector is any data point x_i such that $\frac{|\langle x_i, \beta \rangle|}{\|\beta\|_2} = m$,
where $m = \min_i \frac{|\langle x_i, \beta \rangle|}{\|\beta\|_2}$.

x_i is support vector

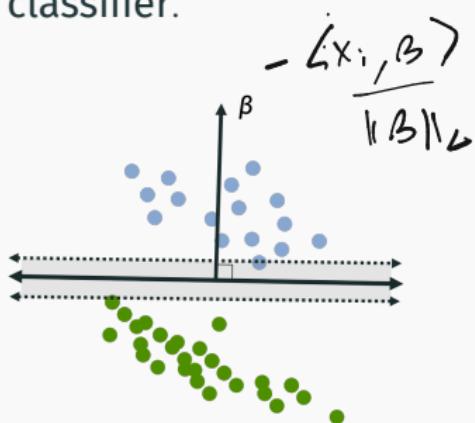
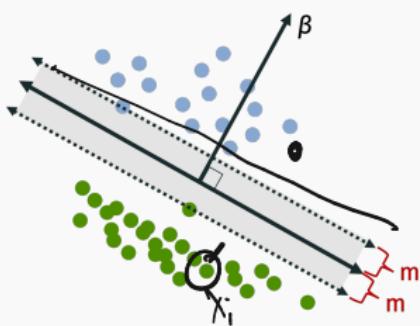
if $\Delta_i = M$

=



HARD-MARGIN SVM

A hard-margin support vector machine (SVM) classifier finds the maximum margin (MM) linear classifier.



i.e. the separating hyperplane which maximizes the margin m . Like regularization, doing so can prevent overfitting even if we can fit our data perfectly.

MARGIN

$$\begin{array}{c} 1/\beta \\ -1 \end{array}$$

$$\min_i \Delta_i = \min_i \frac{|\langle x_i, \beta \rangle|}{\|\beta\|_2}$$

Denote the maximum margin by $\underline{\underline{m^*}}$.

$$\begin{aligned} \underline{\underline{m^*}} &= \max_{\beta} \left[\min_{i \in 1, \dots, n} \frac{|\langle x_i, \beta \rangle|}{\|\beta\|_2} \right] \\ &= \max_{\beta} \left[\min_{i \in 1, \dots, n} \frac{y_i \cdot \langle x_i, \beta \rangle}{\|\beta\|_2} \right] \end{aligned}$$

where $y_i = -1, 1$ depending on what class x_i is in.³

³Note that this is a different convention than the 0, 1 class labels we typically use.

HARD-MARGIN SVM

Original problem: $\max_{\beta} \left[\min_{i \in 1, \dots, n} \frac{y_i \cdot \langle x_i, \beta \rangle}{\|\beta\|_2} \right]$

$\boxed{\langle x_{new}, \beta \rangle > 0}$

Equivalent formulation:

$$\min_{\beta} \|\beta\|_2$$

subject to

$y_i \cdot \langle x_i, \beta \rangle \geq 1$ for all i .

$$\begin{aligned} \max_{\beta} & \left[\min_i \underbrace{y_i \cdot \langle x_i, \beta \rangle}_{\|\beta\|_2} \right] \\ \text{subject to} & y_i \cdot \langle x_i, \beta \rangle \geq 1 \text{ for all } i \\ \text{and} & y_i \cdot \langle x_i, \beta \rangle = 1 \text{ for some } i. \end{aligned}$$

$$= \frac{1}{\|\beta\|_2}$$

Training data:

$$\begin{aligned} \langle x_i, \beta \rangle &\geq 1 \text{ for all } i \text{ s.t. } y_i = 1 \\ \langle x_i, \beta \rangle &\leq -1 \text{ for all } i \text{ s.t. } y_i = -1 \end{aligned}$$

HARD-MARGIN SVM

Equivalent formulation:

$$\min_{\beta} \|\beta\|_2 \quad \text{subject to} \quad y_i \cdot \langle \mathbf{x}_i, \beta \rangle \geq 1 \text{ for all } i.$$

Under this formulation $m^* = \frac{1}{\|\beta\|_2}$.

This is a **constrained optimization problem**. In particular, a linearly constrained quadratic program, which is a type of problem we have efficient optimization algorithms for.

Not as easy to solve as a standard unconstrained, convex optimization problem like logistic regression!

HARD-MARGIN SVM CLASSIFICATION

Classification rule is the same as usual: classify in class 1 if $\langle \mathbf{x}_i, \beta \rangle \geq 0$, class -1 if $\langle \mathbf{x}_i, \beta \rangle < 0$.

Kernel case: As before, we can parameterize as $\beta = \underline{\mathbf{X}^T \alpha}$. When using a non-linear kernel $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$, we have:

$$\underline{\beta = \phi(\mathbf{X})^T \alpha} \quad \langle \phi(\mathbf{x}_{new}), \beta \rangle$$

and to classify a new point $\underline{\mathbf{x}_{new}}$ we compute:

$$\begin{aligned} \underline{\phi(\mathbf{x}_{new})^T \beta} &= \underline{\phi(\mathbf{x}_{new})^T \phi(\mathbf{X})^T \alpha} = \phi(\mathbf{x}_{new})^T \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) \\ &= \sum_{i=1}^n \alpha_i \underline{\phi(\mathbf{x}_{new})^T \phi(\mathbf{x}_i)} \\ &= \boxed{\sum_{i=1}^n \alpha_i k(\mathbf{x}_{new}, \underline{\mathbf{x}_i})}. \end{aligned}$$

Can show that $\alpha_i = 0$ whenever \mathbf{x}_i is not a support vector.

Classification cost scales with # of support vectors, s , not n .

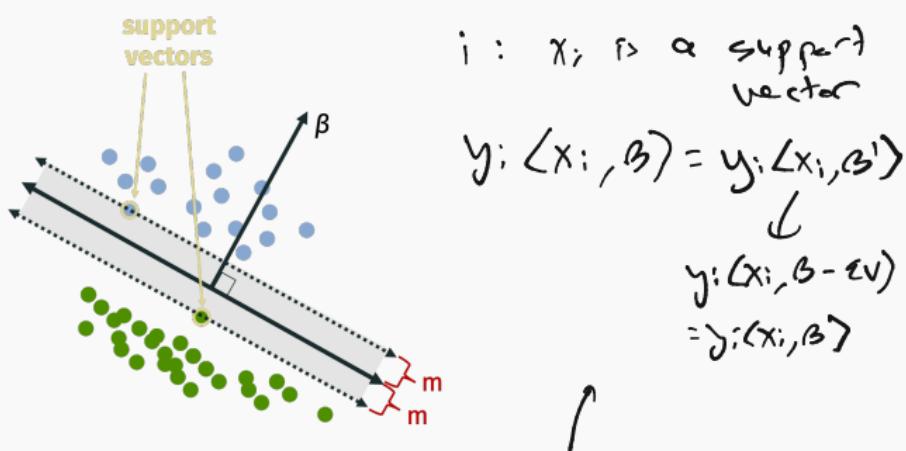
HARD-MARGIN SVM CLASSIFICATION

$$\left(\min_{\beta} \|\beta\|_2 \quad \text{subject to} \quad y_i \cdot \langle x_i, \beta \rangle \geq 1 \text{ for all } i. \right)$$

Claim: Let x_1, \dots, x_s be all vectors with $y_i \cdot \langle x_i, \beta \rangle = 1$. $s < n$

$$\beta = \sum_{i=1}^s \alpha_i x_i.$$

$$\sum_{i \neq s} x_i$$



Proof by contradiction: Suppose $\beta = \sum_{i=1}^s \alpha_i x_i + v$, where v is orthogonal to the support vectors. Can reduced $\|\beta\|_2$ by setting $\beta' \leftarrow \sum_{i=1}^s \alpha_i x_i + (1 - \epsilon)v$ for some small ϵ

HARD-MARGIN SVM CLASSIFICATION

Take-away:

- Training SVMs is typically harder than training using logistic loss.)
- Classification after the model is trained requires $O(n)$ kernel evaluations for general linear classifiers (e.g., found via logistic regression), but just $O(s)$ for an SVM with s support vectors. Often, $s \ll n$.
- Advantages in-terms of storage space as well: only need to keep support vectors around for classification.

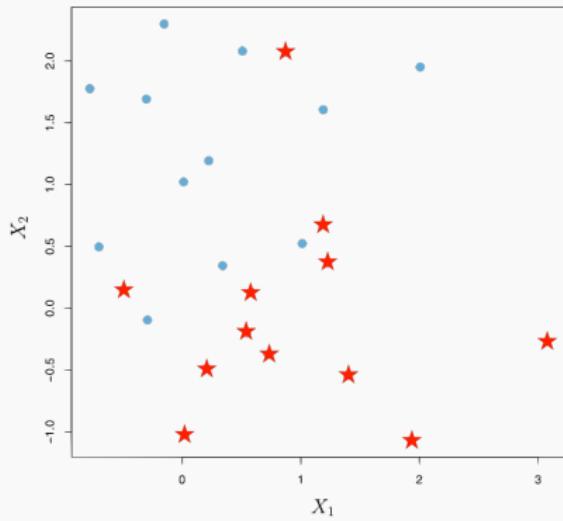
53

Moderate
37 | 53

Mean
38 | 53

HARD-MARGIN SVM

Hard-margin SVMs have a few other critical issues in practice:

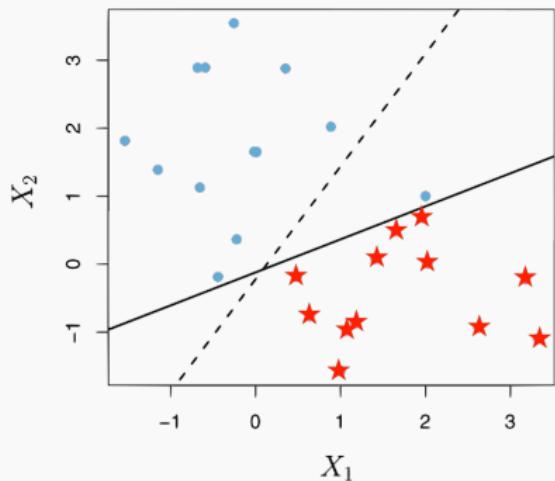
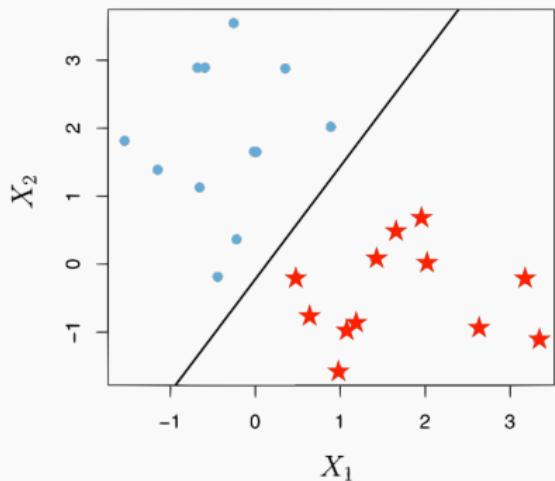


Data might not be linearly separable, in which case the maximum margin classifier is not even defined.

Less likely to be an issue when using a non-linear kernel. If K is full rank then perfect separation is always possible. And typically it is, e.g. for an RBF kernel or moderate degree polynomial kernel.

HARD-MARGIN SVM

Another critical issue in practice:



Hard-margin SVM classifiers are not robust.

SOFT-MARGIN SVM

Solution: Allow the classifier to make some “mistakes”! A mistake can either be a misclassification, or simply a point allowed to be “inside” the margin.

Hard margin objective:

$$\min_{\beta} \|\beta\|_2^2 \quad \text{subject to} \quad y_i \cdot \langle \mathbf{x}_i, \beta \rangle \geq 1 \text{ for all } i.$$

Soft margin objective:

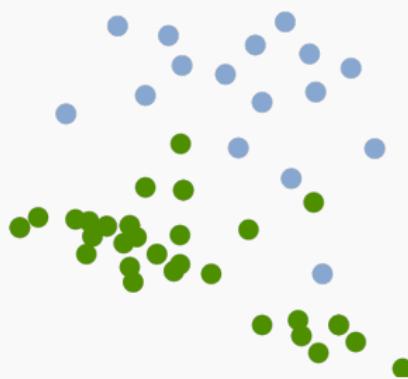
$$\min_{\beta} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad y_i \cdot \langle \mathbf{x}_i, \beta \rangle \geq 1 - \epsilon_i \text{ for all } i.$$

where $\epsilon_i \geq 0$ is a non-negative “slack variable”. This is the magnitude of the “error” (distance past the margin) we allow example \mathbf{x}_i to travel. $\epsilon_i \geq 1$ corresponds to a misclassification.

$C \geq 0$ is a non-negative tuning parameter.

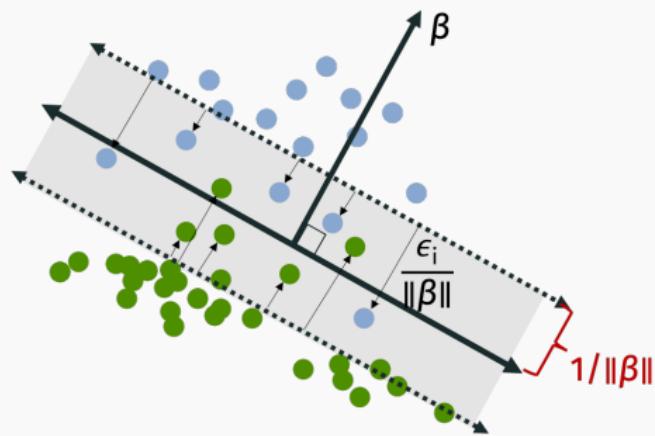
SOFT-MARGIN SVM

Example of a non-separable problem:



SOFT-MARGIN SVM

Recall that $\Delta_i = \frac{y_i \cdot \langle \mathbf{x}_i, \beta \rangle}{\|\beta\|_2}$.

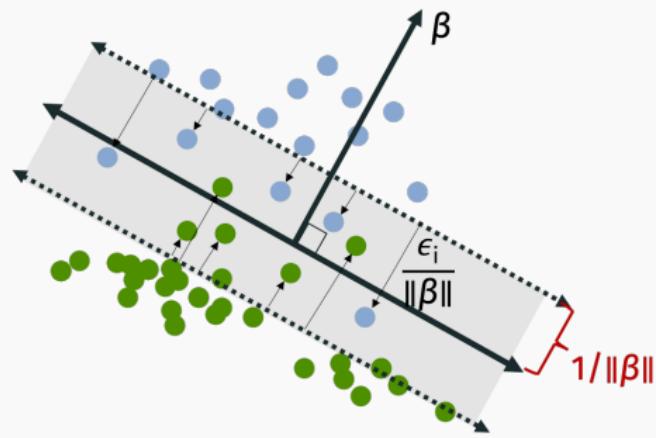


Soft margin objective:

$$\min_{\beta} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad y_i \cdot \langle \mathbf{x}_i, \beta \rangle \geq 1 - \epsilon_i \text{ for all } i.$$

SOFT-MARGIN SVM

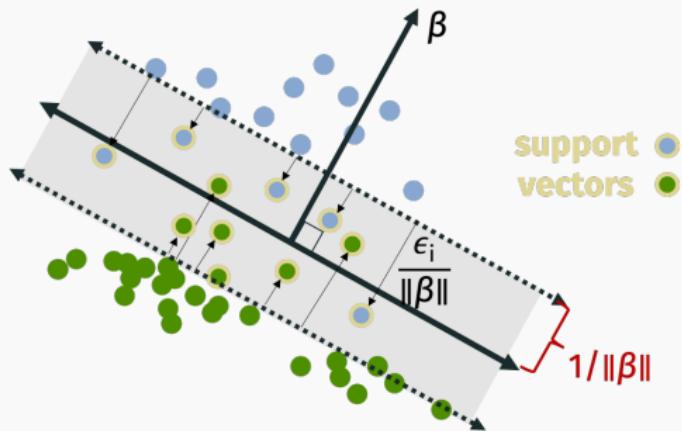
Recall that $\Delta_i = \frac{y_i \cdot \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle}{\|\boldsymbol{\beta}\|_2}$.



Soft margin objective:

$$\min_{\boldsymbol{\beta}} \|\boldsymbol{\beta}\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad \frac{y_i \cdot \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle}{\|\boldsymbol{\beta}\|_2} \geq \frac{1}{\|\boldsymbol{\beta}\|_2} - \frac{\epsilon_i}{\|\boldsymbol{\beta}\|_2} \text{ for all } i.$$

SOFT-MARGIN SVM



Any x_i with a non-zero ϵ_i is a support vector. As before, only support vectors are needed for classification in the kernel setting.

Soft margin objective:

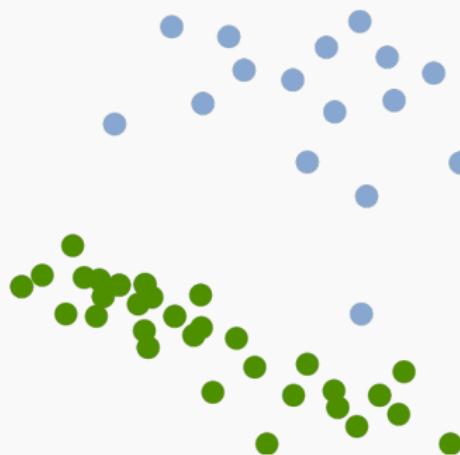
$$\min_{\beta} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i.$$

- Large C means penalties are punished more in objective
 \Rightarrow smaller margin, less support vectors.
- Small C means penalties are punished less in objective
 \Rightarrow larger margin, more support vectors.

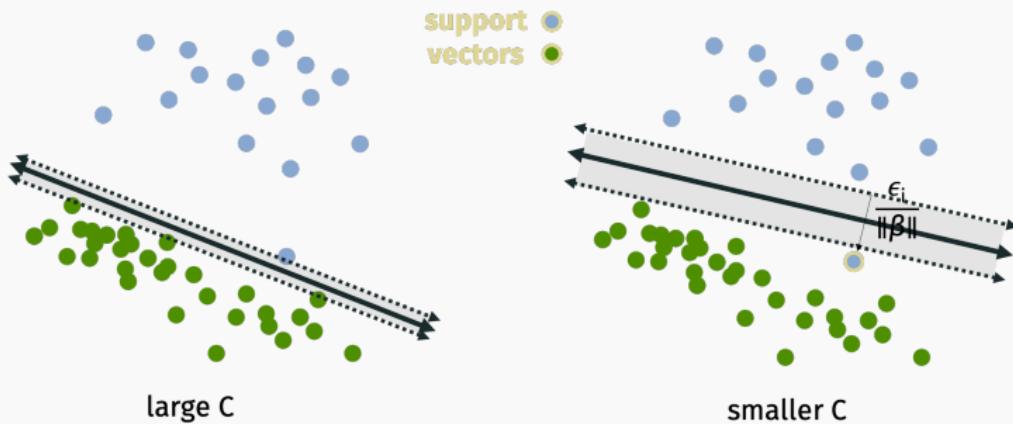
When data is linearly separable, as $C \rightarrow \infty$ we will always get a separating hyperplane. A smaller value of C might lead to a more robust solution.

EFFECT OF C

Example dataset:



EFFECT OF C



The classifier on the right is intuitively more robust. So for this data, a smaller choice for C might make sense.

COMPARISON TO LOGISTIC REGRESSION

Some basic transformations of the soft-margin objective:

$$\min_{\beta} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad y_i \cdot \langle \mathbf{x}_i, \beta \rangle \geq 1 - \epsilon_i \text{ for all } i.$$

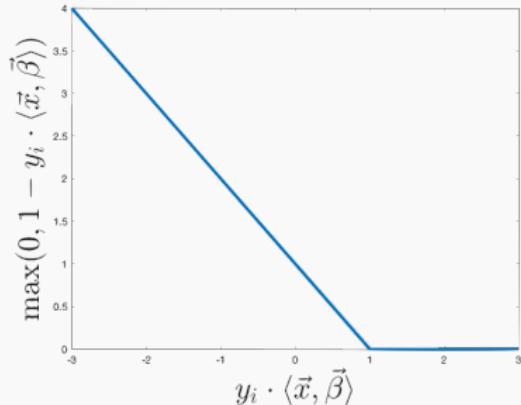
$$\min_{\beta} \|\beta\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i \cdot \langle \mathbf{x}_i, \beta \rangle).$$

$$\min_{\beta} \lambda \|\beta\|_2^2 + \sum_{i=1}^n \max(0, 1 - y_i \cdot \langle \mathbf{x}_i, \beta \rangle).$$

These are all equivalent. $\lambda = 1/C$ is just another scaling parameter.

HINGE LOSS

Hinge-loss: $\max(0, 1 - y_i \cdot \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle)$. Recall that $y_i \in \{-1, 1\}$.



Soft-margin SVM:

$$\min_{\boldsymbol{\beta}} \left[\sum_{i=1}^n \max(0, 1 - y_i \cdot \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle) + \lambda \|\boldsymbol{\beta}\|_2^2 \right]. \quad (1)$$

LOGISTIC LOSS

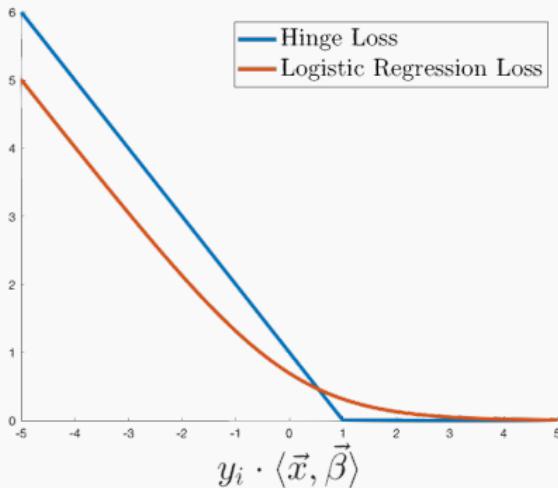
Recall the logistic loss for $y_i \in \{0, 1\}$:

$$\begin{aligned} L(\beta) &= - \sum_{i=1}^n y_i \log(h(\langle \mathbf{x}_i, \beta \rangle)) + (1 - y_i) \log(1 - h(\langle \mathbf{x}_i, \beta \rangle)) \\ &= - \sum_{i=1}^n y_i \log \left(\frac{1}{1 + e^{-\langle \mathbf{x}_i, \beta \rangle}} \right) + (1 - y_i) \log \left(\frac{e^{-\langle \mathbf{x}_i, \beta \rangle}}{1 + e^{-\langle \mathbf{x}_i, \beta \rangle}} \right) \\ &= - \sum_{i=1}^n y_i \log \left(\frac{1}{1 + e^{-\langle \mathbf{x}_i, \beta \rangle}} \right) + (1 - y_i) \log \left(\frac{1}{1 + e^{\langle \mathbf{x}_i, \beta \rangle}} \right) \end{aligned}$$

COMPARISON OF SVM TO LOGISTIC REGRESSION

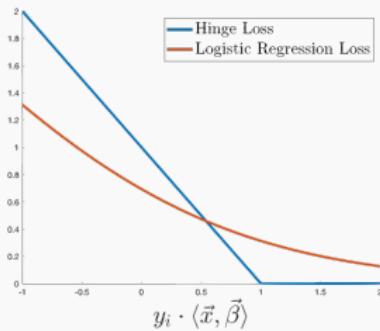
Compare this to the logistic regression loss reformulated for $y_i \in \{-1, 1\}$:

$$\sum_{i=1}^n -\log \left(\frac{1}{1 - e^{-y_i \cdot \langle \vec{x}_i, \vec{\beta} \rangle}} \right)$$



COMPARISON TO LOGISTIC REGRESSION

So, in the end, the function minimized when finding β for the standard **soft-margin SVM** is very similar to the objective function minimized when finding β using **logistic regression with ℓ_2 regularization**. Sort of...



Both functions can be optimized using first-order methods like gradient descent. This is now a common choice for large problems. Will explore more on next lab.