

CS-UY 4563: Lecture 13

Kernel Methods

NYU Tandon School of Engineering, Prof. Christopher Musco

My new office: <https://nyu.zoom.us/my/cmusco>

- Visit this URL during my office hours or any individual meetings.
- You can also “drop in” even if I’m not in the chat. I will receive an email notification (might take me a few minutes to notice) and can then let you in if I’m free.
- For lectures still use the links on NYU Classes (which allows for automatic recording, transcription, etc.)

By 4/1 (next Wednesday) choose a partner and topic.

- Email me team members, project topic, and a sentence or two about your idea. If you have any data sets in mind, let me know that as well.
- Then set up a meeting at: https://docs.google.com/spreadsheets/d/1DsR7ia4VfYb5joIavsG8_T1JBgAufkVbdT7bLfPGGQ0/edit?usp=sharing.

LETS EASE BACK INTO THINGS

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

Training data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ where $y_1, \dots, y_n \in \{1, \dots, q\}$.

Classification algorithm:

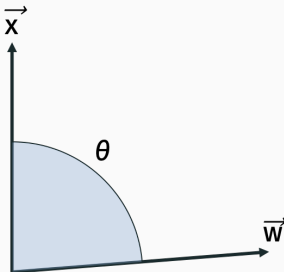
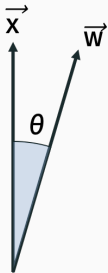
Given new input \vec{x}_{new} ,

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

INNER PRODUCT SIMILARITY

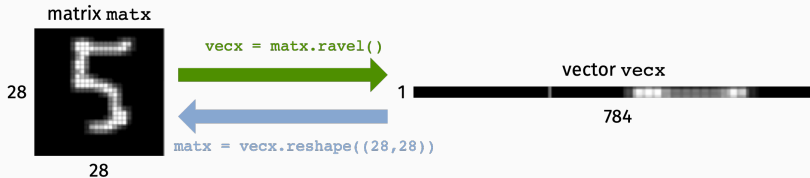
Given data vectors $\vec{x}, \vec{w} \in \mathbb{R}^d$, the inner product $\langle \vec{x}, \vec{w} \rangle$ is a natural similarity measure.

$$\langle \vec{x}, \vec{w} \rangle = \sum_{i=1}^d \vec{x}[i] \vec{w}[i] = \cos(\theta) \|\vec{x}\|_2 \|\vec{w}\|_2.$$



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.

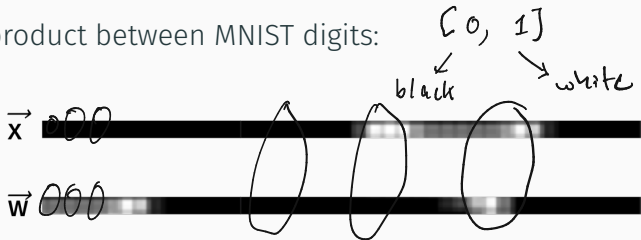


```
xmat = np.array([[1,2,3],[4,5,6],[7,8,9]])  
  
array([[1, 2, 3],  
       [4, 5, 6],  
       [7, 8, 9]])
```

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

INNER PRODUCT FOR MNIST

Inner product between MNIST digits:

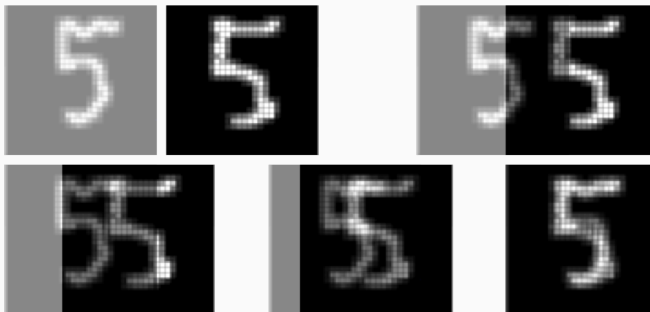


$$\langle \vec{x}, \vec{w} \rangle = \sum_{i=1}^{28} \sum_{j=1}^{28} \underline{\text{mat}x[i,j] \text{mat}w[i,j]}.$$

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.

INNER PRODUCT FOR MNIST

Visualizing the inner product between two images:



Images with high inner product have a lot of overlap.

cf 10

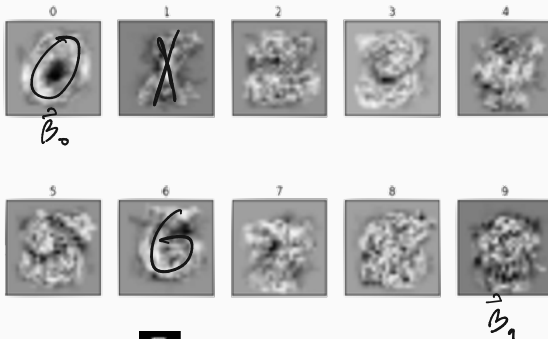
One-vs.-all Classification with Logistic Regression:

- Learn q classifiers with parameters $\vec{\beta}_0, \vec{\beta}_1, \dots, \vec{\beta}_{q-1}$.
- Given \vec{x}_{new} compute $\langle \vec{x}_{new}, \vec{\beta}_0 \rangle, \dots, \langle \vec{x}_{new}, \vec{\beta}_{q-1} \rangle$.
- Predict class $y_{new} = \arg \max_i \langle \vec{x}_{new}, \vec{\beta}_i \rangle$.

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

MATCHED FILTER

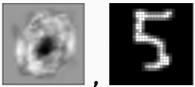
Visualizing $\vec{\beta}_0, \dots, \vec{\beta}_9$:

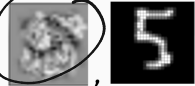


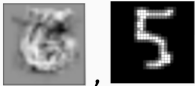
For an input image **5**, compute inner product similarity with all weight matrices and choose most similar one.

In contrast to k -NN, only need to compute similarity with q items instead of n .

VIEW OF LOGISTIC REGRESSION

$$\langle \vec{\beta}_0, \vec{x}_{new} \rangle = 45$$
The image shows a grayscale digit '0' in a square frame, followed by a comma, a black square frame containing a white digit '5', and then the text '= 45'. The digit '0' is slightly blurry.

$$\langle \vec{\beta}_5, \vec{x}_{new} \rangle = 212$$
The image shows a grayscale digit '5' in a square frame, followed by a comma, a black square frame containing a white digit '5', and then the text '= 212'. The digit '5' in the first frame is circled in black.

$$\langle \vec{\beta}_6, \vec{x}_{new} \rangle = 84$$
The image shows a grayscale digit '5' in a square frame, followed by a comma, a black square frame containing a white digit '5', and then the text '= 84'. The digit '5' in the first frame is slightly blurry.

Select class for \vec{x}_{new} which achieves highest “score”, as measured by the inner product similarity.

Logistic Regression Model:

Given data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ (here $d = 784$) and binary label vector $\vec{y} \in \{0, 1\}^n$ for class i (1 if in class i , 0 if not), find $\vec{\beta}_i \in \mathbb{R}^d$ to minimize the log loss between:

$$\underline{\vec{y}}$$

and

$$h(\mathbf{X}\underline{\vec{\beta}})$$

↗ data matrix

where $h(\mathbf{X}\vec{\beta}_i) = \frac{1}{1+e^{-\mathbf{X}\vec{\beta}_i}}$ applies the logistic function entrywise to the vector $\mathbf{X}\vec{\beta}$.

$$\text{Loss} = - \sum_{j=1}^n y_j \log(h(\mathbf{X}\vec{\beta}_i)[j]) + (1 - y_j) \log(1 - h(\mathbf{X}\vec{\beta}_i)[j])$$

Logistic Regression Model:

Given data matrix $X \in \mathbb{R}^{n \times d}$ (here $d = 784$) and binary label vector $\vec{y} \in \{0, 1\}^n$ for class i (1 if in class i , 0 if not), find $\vec{\beta} \in \mathbb{R}^d$ to minimize the log loss between:

$$\underline{\underline{\vec{y}}}$$

and

$h(X\vec{\beta})$
logistic function

Reminder from linear algebra: Without loss of generality, can assume that $\vec{\beta}$ lies in the row span of X .

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

$$\vec{\beta} = X^T \vec{\alpha}.$$

$$X\beta = XX^T\alpha$$

Logistic Regression Equivalent Formulation:

Given data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ (here $d = 784$) and binary label vector $\vec{y} \in \{0, 1\}^n$ for class i (1 if in class i , 0 if not), find $\vec{\alpha} \in \mathbb{R}^n$ to minimize the log loss between:

\vec{y}
 $\{0, 0, 1, 1, 0, 1, \dots\}$

and

$h(\mathbf{X}\mathbf{X}^T\vec{\alpha})$.

Can still be minimized via gradient descent:

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

What does classification for a new point \vec{x}_{new} look like? Recall that for a given one-vs-all classification i , the original parameter vector $\underline{\vec{\beta}_i} = \underline{\mathbf{X}^T \vec{\alpha}_i}$ $\vec{\alpha}_0, \dots, \vec{\alpha}_q$

- Learn q classifiers with parameters $\underline{\vec{\alpha}_1, \vec{\alpha}_2, \dots, \vec{\alpha}_q}$.
- Given \vec{x}_{new} compute $\langle \vec{x}_{new}, \underline{\mathbf{X}^T \vec{\alpha}_0} \rangle, \dots, \langle \vec{x}_{new}, \mathbf{X}^T \vec{\alpha}_q \rangle$,
- Predict class $y_{new} = \arg \max_i \langle \vec{x}_{new}, \mathbf{X}^T \vec{\alpha}_i \rangle$.

REFORMULATED VIEW

$$\langle \vec{x}_{new}, \vec{b}_i \rangle$$

\downarrow \downarrow
 (727×1) $(727, 1)$

Score for class i :

$$X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{727}$$

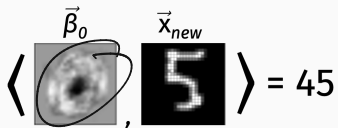
$$\begin{aligned}
 \langle \vec{x}_{new}, X^T \vec{\alpha}_i \rangle &= (\vec{x}_{new}^T X) \vec{\alpha}_i \\
 &= \langle X \vec{x}_{new}, \vec{\alpha}_i \rangle \\
 &= \sum_{j=1}^n \alpha_i[j] \langle \vec{x}_{new}, \vec{x}_j \rangle.
 \end{aligned}$$

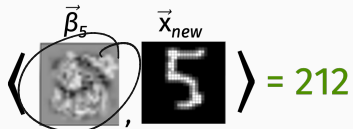
$\underbrace{\hspace{10em}}_{727 \times 1}$

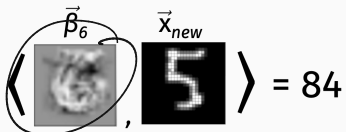
$$\begin{aligned}
 (727 \times n) \cdot (n \times 1) \\
 = (727 \times 1)
 \end{aligned}$$

weight for every \vec{x}_j in training data







ORIGINAL VIEW OF LOGISTIC REGRESSION

$$\langle \vec{\beta}_0, \vec{x}_{new} \rangle = 45$$


$$\langle \vec{\beta}_5, \vec{x}_{new} \rangle = 212$$


$$\langle \vec{\beta}_6, \vec{x}_{new} \rangle = 84$$


NEW VIEW OF LOGISTIC REGRESSION

| | \vec{x}_1 | \vec{x}_2 | \vec{x}_3 | \vec{x}_4 | \vec{x}_5 | \vec{x}_6 | ... |
|------------------|---|---|---|---|---|---|-----|
| |  |  |  |  |  |  | ... |
| $\vec{\alpha}_0$ | $\vec{\alpha}_0[1] = .4$ | $\vec{\alpha}_0[2] = .9$ | $\vec{\alpha}_0[3] = .2$ | $\vec{\alpha}_0[4] = .1$ | $\vec{\alpha}_0[5] = .8$ | $\vec{\alpha}_0[6] = .05$ | ... |
| $\vec{\alpha}_1$ | $\vec{\alpha}_1[1] = .05$ | $\vec{\alpha}_1[2] = .1$ | $\vec{\alpha}_1[3] = .2$ | $\vec{\alpha}_1[4] = .02$ | $\vec{\alpha}_1[5] = .1$ | $\vec{\alpha}_1[6] = .95$ | ... |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | |
| $\vec{\alpha}_5$ | $\vec{\alpha}_5[1] = .1$ | $\vec{\alpha}_5[2] = .2$ | $\vec{\alpha}_5[3] = .85$ | $\vec{\alpha}_5[4] = .75$ | $\vec{\alpha}_5[5] = .1$ | $\vec{\alpha}_5[6] = .05$ | ... |
| | \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | |

Learn n length parameter vectors $\vec{\alpha}_0, \dots, \vec{\alpha}_9$, one for each class.

NEW VIEW OF LOGISTIC REGRESSION

$$\sum_{j=0}^M \vec{\sigma}_j[j] \langle \vec{x}_j, \vec{x}_{new} \rangle \quad \text{for } i = 0, \dots, 9$$

$$\vec{\alpha}_0[1] \times \langle \vec{x}_1, \vec{x}_{new} \rangle + \vec{\alpha}_0[2] \times \langle \vec{x}_2, \vec{x}_{new} \rangle + \vec{\alpha}_0[3] \times \langle \vec{x}_3, \vec{x}_{new} \rangle + \dots = 45$$

$$\vec{\alpha}_5[1] \times \langle \vec{x}_1, \vec{x}_{new} \rangle + \vec{\alpha}_5[2] \times \langle \vec{x}_2, \vec{x}_{new} \rangle + \vec{\alpha}_5[3] \times \langle \vec{x}_3, \vec{x}_{new} \rangle + \dots = 212$$

$$\vec{\alpha}_6[1] \times \langle \vec{x}_1, \vec{x}_{new} \rangle + \vec{\alpha}_6[2] \times \langle \vec{x}_2, \vec{x}_{new} \rangle + \vec{\alpha}_6[3] \times \langle \vec{x}_3, \vec{x}_{new} \rangle + \dots = 84$$

Classification looks similar to k -NN: we compute the similarity between \vec{x}_{new} and every other vector in our training data set. A weighted sum of the similarities leads to scores for each class.

Assign \vec{x}_{new} to the class with highest score.

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):

$\sum_{j=1}^{724} z_j x_j$
 $z_j \leq 1$

$$\langle \vec{z}, \vec{x} \rangle < \langle \vec{y}, \vec{x} \rangle$$

But clearly the first image is more similar.

$$\langle \vec{z}, \vec{x} \rangle < \langle \vec{y}, \vec{x} \rangle$$

Here's a more realistic scenario.

KERNEL FUNCTIONS: PERSPECTIVE ONE

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

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

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

Gaussian kernel:

$$e^{-\|\vec{x}-\vec{y}\|_2^2} < e^{-\|\vec{x}-\vec{z}\|_2^2} \quad \text{closer to 0} \quad \text{closer to 1} \quad k(\vec{x}, \vec{y}) = e^{-\frac{\|\vec{x}-\vec{y}\|_2^2}{\sigma^2}} \quad \text{some scaling factor } \sigma^2$$

for some scaling factor σ .

$$\|\vec{x}-\vec{y}\|_2^2 > \|\vec{x}-\vec{z}\|_2^2$$

is large

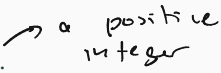
$$k\left(\begin{array}{c} \vec{z} \\ \text{5} \end{array}, \begin{array}{c} \vec{x} \\ \text{5} \end{array}\right) > k\left(\begin{array}{c} \vec{y} \\ \text{5} \end{array}, \begin{array}{c} \vec{x} \\ \text{5} \end{array}\right)$$

$\vec{y} - \vec{x} = \vec{x} - \vec{y}$

→ 0 when $\|\vec{x}-\vec{y}\|_2^2$ is large

→ 1 when $\|\vec{x}-\vec{y}\|_2^2$ is small.

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

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

But you can imagine much more complex similarity metrics.

KERNEL FUNCTIONS: PERSPECTIVE TWO

For a simple algorithm like k -NN you can swap out the inner product similarity with any similarity function you could possibly imagine.

For a methods like logistic regression, this is not the case...

Recall: We learned a parameter vector $\vec{\alpha}$ to minimize $LL(\vec{y}, \mathbf{X}^T \vec{\alpha})$ where $LL()$ denotes the logistic loss. Then we classified via:

$$\underline{\underline{\chi \hat{b}}} = \chi(\mathbf{x}^T \vec{\alpha})$$

$$\langle \vec{x}_{new}, \mathbf{X}^T \vec{\alpha} \rangle = \vec{x}_{new}^T \mathbf{X}^T \vec{\alpha} = \sum_{j=1}^n \underline{\vec{\alpha}[j]} \langle \underline{\vec{x}_{new}}, \underline{\vec{x}_j} \rangle.$$

The inner product similarity came from the fact that our predictions were based on the linear function $\mathbf{X}^T \vec{\alpha}$.

KERNEL FUNCTIONS AS FEATURE TRANSFORMATION

$$\begin{array}{c} \boxed{} \\ x \end{array} \xrightarrow{\phi} \begin{array}{c} \boxed{} \\ \phi(x) \end{array}$$

$$\begin{array}{c} \boxed{} \\ x \end{array} \xrightarrow{\Phi} \begin{array}{c} \boxed{} \\ \Phi(x) \end{array}$$

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

$$\underline{k(\vec{x}, \vec{w})} = \phi(\vec{x})^T \phi(\vec{w})$$

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

$$\begin{array}{c} \boxed{} \\ x \end{array} \xrightarrow{\phi} \begin{array}{c} \boxed{} \\ \underline{\phi(x)} \end{array}$$

$$\begin{array}{c} \boxed{} \\ \omega \end{array} \xrightarrow{\Phi} \begin{array}{c} \boxed{} \\ \underline{\underline{\Phi(\omega)}} \end{array}$$

KERNEL FUNCTIONS AND FEATURE TRANSFORMATION

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

$$d = 3$$

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

$$\phi(\vec{x}) =$$

$$\begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_3 \\ x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \sqrt{2}x_2x_3 \end{bmatrix}$$

$$u = 10$$

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

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

$$+ 2x_1 w_1 x_2 w_2 + 2x_1 w_1 x_3 w_3 + 2x_2 w_2 x_3 w_3$$

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

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

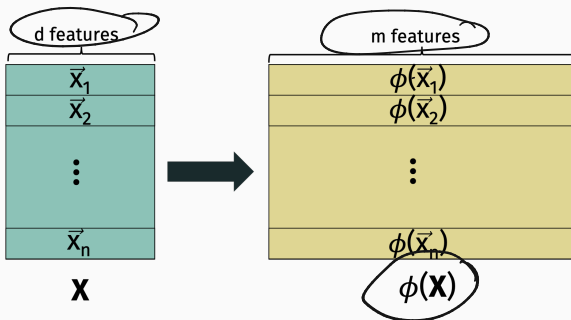
- Gaussian RBF Kernel: $k(\vec{x}, \vec{y}) = e^{-\|\vec{x} - \vec{y}\|_2^2 / \sigma^2}$
- Laplace Kernel: $k(\vec{x}, \vec{y}) = e^{-\|\vec{x} - \vec{y}\|_2 / \sigma}$
- Polynomial Kernel: $k(\vec{x}, \vec{y}) = (\langle \vec{x}, \vec{y} \rangle + 1)^q$

And there are many more...

KERNEL FUNCTIONS AND FEATURE TRANSFORMATION

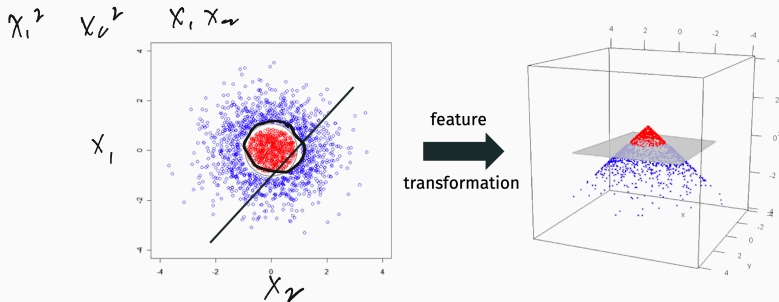
Feature transformations \iff new similarity metrics.

Using $k(\cdot, \cdot)$ in place of the inner product $\langle \cdot, \cdot \rangle$ is **equivalent** to replacing every data point $\vec{x}_1, \dots, \vec{x}_n$ in our data set with $\phi(\vec{x}_1), \dots, \phi(\vec{x}_n)$.



TAKEAWAY ONE

We can improve performance by replacing the inner product with another kernel $k(\cdot, \cdot)$ for the same reason that feature transformations improved performance.



When you add features, it becomes possible to learn more complex decision boundaries (in this case a circle) with a linear classifier.

PSD kernel functions give a principled way of “swapping out” the inner product with a new similarity metric for linear algorithms like multiple linear regression or logistic regression. For non-PSD kernels it is not clear how to do this.

KERNEL LOGISTIC REGRESSION

Standard logistic regression

Loss function: $\underline{XX^T}$

$$L(\vec{\alpha}) = \overset{\text{Loss}}{LL}(\vec{y}, \underline{XX^T} \vec{\alpha}).$$

Gradient:

$$\nabla L(\vec{\alpha}) = \underline{XX^T} (h(\underline{XX^T} \vec{\alpha}) - \vec{y}).$$

Prediction:

$$\langle \vec{\beta}, \vec{x}_{new} \rangle$$
$$z = \sum_{j=1}^n \vec{\alpha}[j] \langle \vec{x}_{new}, \vec{x}_j \rangle.$$

$$y_{new} = \mathbb{1}[z > 0]$$

Kernel logistic regression

$$\phi(X) \phi(X)^T$$

Loss function:

$$L(\vec{\alpha}) = LL(\vec{y}, \overset{\phi(X)}{\phi(X)^T} \vec{\alpha}).$$

Gradient:

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

Prediction:

$$z = \sum_{j=1}^n \vec{\alpha}[j] \langle \phi(\vec{x}_{new}), \phi(\vec{x}_j) \rangle$$

$\parallel K(\vec{x}_{new}, \vec{x}_j)$

$$y_{new} = \mathbb{1}[z > 0]$$

$$e^{-\|\vec{x}_{new} - \vec{x}_j\|^2}$$

KERNEL REGRESSION

Standard linear regression

Loss function: $\|\vec{y} - \mathbf{X}\vec{\alpha}\|_2$

$$L(\vec{\alpha}) = \|\vec{y} - \mathbf{X}\mathbf{X}^T\vec{\alpha}\|_2$$

Gradient: $\hat{\mathbf{X}}^T \alpha = \vec{b}$

$$\nabla L(\vec{\alpha}) = \underline{2\mathbf{X}\mathbf{X}^T(\mathbf{X}\vec{\alpha} - \vec{y})}.$$

Prediction:

$$y_{new} = \sum_{j=1}^n \vec{\alpha}[j] \langle \vec{x}_{new}, \vec{x}_j \rangle.$$

$$\int : \mathbf{X}\mathbf{X}^T \vec{v}$$

Kernel linear regression

Loss function:

$$L(\vec{\alpha}) = \|\vec{y} - \phi(\mathbf{X})\phi(\mathbf{X})^T\vec{\alpha}\|_2$$

Gradient:

$$\nabla L(\vec{\alpha}) = \underline{2\phi(\mathbf{X})\phi(\mathbf{X})^T(\phi(\mathbf{X})\phi(\mathbf{X})^T\vec{\alpha} - \vec{y})}.$$

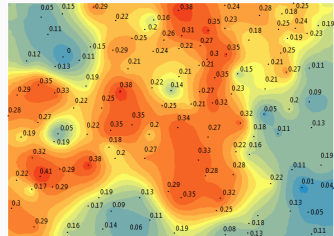
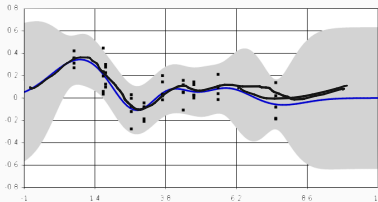
Prediction:

$$y_{new} = \sum_{j=1}^n \vec{\alpha}[j] \langle \phi(\vec{x}_{new}), \phi(\vec{x}_j) \rangle.$$

$\nearrow \kappa(\vec{x}_{new}, \vec{x}_j)$

KERNEL REGRESSION

We won't study kernel regression in detail, but it's a very important statistical tool, especially when dealing with spatial or temporal data.

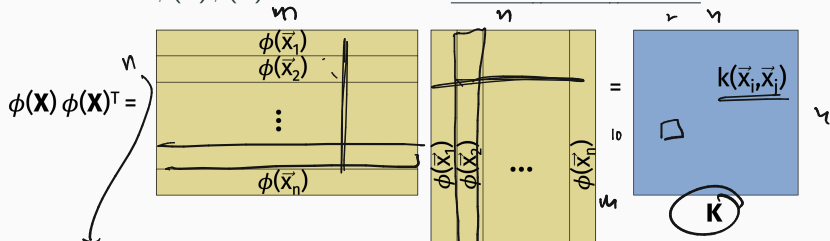


Also known as Gaussian Process (GP) Regression or Kriging.

KERNEL MATRIX

$X X^T$ "Gram matrix"
or "Gramian matrix"

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



number of
training data points



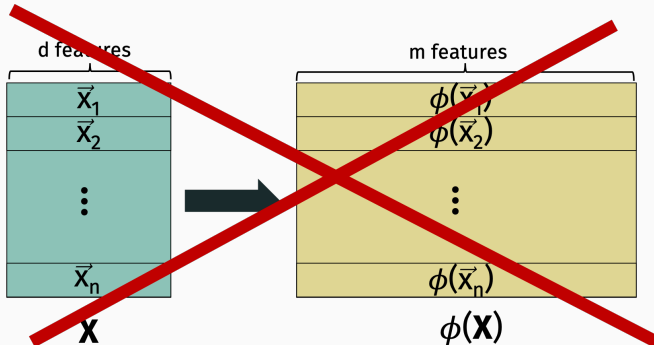
$$K(\vec{x}_2, \vec{x}_{10})$$

$$= \langle \phi(\vec{x}_2), \phi(\vec{x}_{10}) \rangle$$

KERNEL TRICK

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

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



This can lead to significant computational savings!

- 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*.

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

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