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.

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

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 between MNIST digits:



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

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.

Visualizing the inner product between two images:



Images with high inner product have a lot of overlap.

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 × 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 $\vec{\beta}_0, \ldots, \vec{\beta}_9$:





For an input image **5**, compute <u>inner product</u> 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



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:

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

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 $\mathbf{X} \in \mathbb{R}^{n \times d}$ (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 $\vec{\beta} \in \mathbb{R}^d$ to minimize the log loss between:



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

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

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

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:

$$\vec{y}$$
 and $h(XX^T\vec{\alpha})$.

Can still be minimized via gradient descent:

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

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

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

Score for class *i*:

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

ORIGINAL VIEW OF LOGISTIC REGRESSION



NEW VIEW OF LOGISTIC REGRESSION



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

NEW VIEW OF LOGISTIC REGRESSION



Classification looks similar to *k*-NN: we compute the <u>similarity</u> 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.

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

$$\langle \mathbf{5}, \mathbf{5} \rangle < \langle \mathbf{5}, \mathbf{5} \rangle$$

But clearly the first image is more similar.

$$\langle \begin{array}{c} \vec{z} \\ \vec{q} \\ \vec{q} \\ \vec{q} \\ \end{pmatrix} < \langle \begin{array}{c} \vec{y} \\ \vec{y} \\ \vec{y} \\ \vec{q} \\ \end{pmatrix} \rangle$$

Here's a more realistic scenario.

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

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

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

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

for some scaling factor σ .

$$k(\overbrace{5}^{\vec{z}},\overbrace{5}^{\vec{x}}) > k([,\overbrace{5}^{\vec{y}},\overbrace{5}^{\vec{x}}])$$

Lots of kernel functions 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: $k(\vec{x}, \vec{y}) = (\langle \vec{x}, \vec{y} \rangle + 1)^q$.

But you can imagine much more complex similarity metrics.

For a simple algorithm like *k*-NN you can swap our the inner product similarity with <u>any similarity function you could</u> <u>possibly imagine</u>.

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:

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

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

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

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

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

KERNEL FUNCTIONS AND FEATURE TRANSFORMATION

Example: Degree 2 polynomial kernel, $k(\vec{x}, \vec{w}) = (\vec{x}^T \vec{w} + 1)^2$. $\phi(\vec{x}) = \begin{bmatrix} 1\\ \sqrt{2}x_1\\ \sqrt{2}x_2\\ \sqrt{2}x_3\\ x_1^2\\ x_2^2\\ x_3^2\\ \sqrt{2}x_1x_2\\ \sqrt{2}x_1x_3\\ \sqrt{2}x_2x_3 \end{bmatrix}$ $\vec{X} = \begin{bmatrix} X_1 \\ X_2 \\ X_2 \end{bmatrix}$

$$(\vec{x}^T \vec{w} + 1)^2 = (x_1 y_1 + x_2 y_2 + x_3 y_3 + 1)^2$$

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

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 ...

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 $\vec{x}_1, \ldots, \vec{x}_n$ in our data set with $\phi(\vec{x}_1), \ldots, \phi(\vec{x}_n)$.



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 logisitic regression

Kernel logisitic regression

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

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

Loss function:

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

Gradient:

Gradient:

Loss function:

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

Prediction:

Prediction:

$$z = \sum_{j=1}^{n} \vec{\alpha}[j] \langle \vec{x}_{new}, \vec{x}_{j} \rangle. \qquad z = \sum_{j=1}^{n} \vec{\alpha}[j] \langle \phi(\vec{x}_{new}), \phi(\vec{x}_{j}) \rangle$$
$$y_{new} = \mathbb{1}[z > 0] \qquad y_{new} = \mathbb{1}[z > 0]$$

KERNEL REGRESSION

Standard linear regression

Loss function:

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

Gradient:

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

Prediction:

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

Kernel linear regression

Loss function:

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

Gradient:

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

Prediction:

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

We won't study kernel <u>regression</u> 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.

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



KERNEL TRICK

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

- For training we just need the kernel matrix **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 **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.)