PHY 835: Collider Physics Phenomenology

Machine Learning in Fundamental Physics

Gary Shiu, UW-Madison



Lecture 9: Kernel Methods

Recap of Lecture 8

- Support Vector Machines: functional and geometric margins
- A simple example and animation:

https://www.youtube.com/watch?v=5zRmhOUjjGY

- Optimal margin classifier
- Lagrange duality
- Kernel Methods

Outline for today

- Kernel methods
- Soft margins
- Ensemble Methods
 - Bagging
 - Boosting
 - Random Forest

Ref: Andrew Ng's Lecture Notes: <u>https://sgfin.github.io/files/notes/</u> <u>CS229_Lecture_Notes.pdf</u> and 1803.08823

Feature Maps

• So far we have been considering linear separable data, what if the output is more accurately represented by a non-linear function e.g.

$$y = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0$$

• Consider the function $\phi : \mathbb{R} \to \mathbb{R}^4$ known as the feature map:



• y is a linear function over $\phi(x)$: $y = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0 = \theta^T \phi(x)$.

• For linear functions, our algorithm for the optimal margin classifier requires computing $\langle x, z \rangle$. This suggests that we should:

 $\langle x, z \rangle \rightarrow K(x, z) = \langle \phi(x), \phi(z) \rangle = \phi(x)^T \phi(z)$

- The kernel K(x, z) may be inexpensive to compute, even though it is expensive to calculate $\phi(x)$ (extremely high dimensional vector).
- Using K(x, z), we can get SVM to learn high dimensional feature space given by ϕ without having to explicitly find/represent $\phi(x)$.
- Suppose $x, z \in \mathbb{R}^n$ and the Kernel is given by:

$$K(x,z) = (x^T z)^2.$$

• What is the corresponding feature map $\phi(x)$?

• We can write the Kernel as:

$$K(x,z) = \left(\sum_{i=1}^{n} x_i z_i\right) \left(\sum_{j=1}^{n} x_j z_i\right) = \sum_{i=1}^{n} \sum_{j=1}^{n} x_i x_j z_i z_j = \sum_{i,j=1}^{n} (x_i x_j) (z_i z_j)$$

• The feature map (for n=3):

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

- Computing $\phi(x)$ takes $\mathcal{O}(n^2)$ time while computing K(x, z) takes only $\mathcal{O}(n)$ time (linear in the dim. of the input attributes).
- For a huge input dimension (say \sim 1000), this is a significant speedup.

• For a related Kernel:

$$K(x,z) = (x^{T}z + c)^{2}$$

= $\sum_{i,j=1}^{n} (x_{i}x_{j})(z_{i}z_{j}) + \sum_{i=1}^{n} (\sqrt{2c}x_{i})(\sqrt{2c}z_{i}) + c^{2}.$
$$\begin{bmatrix} x_{1}x_{1} \\ x_{1}x_{2} \\ x_{1}x_{3} \end{bmatrix}$$

• The corresponding feature map (for n=3):

More generally $K(x, z) = (x^T z + c)^d$:

Entries of ϕ = monomials of the form $x_{i_1}x_{i_2}...x_{i_d}$ up to degree d

• Computing K(x, z) still takes only $\mathcal{O}(n)$ time, even for more general:

$$K(x,z) = (x^T z + c)^d$$

	$x_{1}x_{3}$
	$x_2 x_1$
	$x_2 x_2$
	$x_{2}x_{3}$
$\phi(x) =$	$x_{3}x_{1}$
	$x_{3}x_{2}$
	$x_{3}x_{3}$
	$\sqrt{2cx_1}$
	$\sqrt{2cx_2}$
	$\sqrt{2c}x_3$
	c

- Intuitively, if $\phi(x)$ and $\phi(z)$ are close together, $K(x, z) = \phi(x)^T \phi(z)$ is large; conversely if they are far apart, K(x, z) is small.
- K(x, z) is a measure of how similar x and z are; consider e.g.,

$$K(x,z) = \exp\left(-\frac{||x-z||^2}{2\sigma^2}\right)$$

- Can this be a kernel? Yes, Gaussian kernel.
- In general, can we tell if some function K(x, z) is a valid kernel?
- Given a finite set of *m* points, define an $m \times m$ Kernel matrix:

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

• The kernel matrix should be symmetric, and positive semi-definite.

• $K_{ij} = K_{ji}$ is straightforward whereas the positivity of the *K* matrix:

$$z^{T}Kz = \sum_{i} \sum_{j} z_{i}K_{ij}z_{j}$$

$$= \sum_{i} \sum_{j} z_{i}\phi(x^{(i)})^{T}\phi(x^{(j)})z_{j}$$

$$= \sum_{i} \sum_{j} z_{i} \sum_{k} \phi_{k}(x^{(i)})\phi_{k}(x^{(j)})z_{j}$$

$$= \sum_{k} \sum_{i} \sum_{j} z_{i}\phi_{k}(x^{(i)})\phi_{k}(x^{(j)})z_{j}$$

$$= \sum_{k} \left(\sum_{i} z_{i}\phi_{k}(x^{(i)})\right)^{2}$$

$$\geq 0.$$

- Since z is arbitrary, the kernel matrix K is positive semi-definite.
- These two conditions turn out to be not only necessary but sufficient.

Mercer's Theorem

Theorem (Mercer). Let $K : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$ be given. Then for K to be a valid (Mercer) kernel, it is necessary and sufficient that for any $\{x^{(1)}, \ldots, x^{(m)}\}, (m < \infty)$, the corresponding kernel matrix is symmetric positive semi-definite.

- Check the validity of a kernel without constructing the feature map.
- The kernel trick has wider applications than SVM, e.g.
 - In classifying digits (MNIST database), a simple polynomial kernel or the Gaussian kernel gives extremely good performance.
 - In classifying strings of English alphabets with length k, we have a 26^{k} dim. feature space. Kernel calculations take $\mathcal{O}(26)$ time.
- Kernel trick: replacing $\langle x, z \rangle$ by K(x, z) turns the algorithm to work in higher dim. feature space.

SVM: Soft Margins

- Earlier discussion (w/o kernel) assumed data is linearly separable.
- Mapping data to a higher dim. feature space generally increases the chance for the data to be separable, but there is no guarantee.
- Also, choice of separating hyperplane is susceptible to outliners.



SVM: Soft Margins

• Optimization with a soft margin (and penalty):

$$\min_{\gamma,w,b} \quad \frac{1}{2} ||w||^2 + C \sum_{i=1}^m \xi_i$$

s.t. $y^{(i)}(w^T x^{(i)} + b) \ge 1 - \xi_i, \quad i = 1, \dots, m$
 $\xi_i \ge 0, \quad i = 1, \dots, m.$

- Here we introduced an L_1 regularization; we allow the margin to be less than 1 by paying a cost controlled by C.
- This SVM implementation is in the package sklearn.
- We can construct the Lagrangian:

0 0

$$\mathcal{L}(w,b,\xi,\alpha,r) = \frac{1}{2}w^T w + C\sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y^{(i)}(x^T w + b) - 1 + \xi_i \right] - \sum_{i=1}^m r_i \xi_i.$$

• Like before, we can turn this primal problem into its dual problem.

SVM: Soft Margins

• The dual optimization problem:

$$\max_{\alpha} \quad W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

s.t. $0 \le \alpha_i \le C, \quad i = 1, \dots, m$
 $\sum_{i=1}^{m} \alpha_i y^{(i)} = 0,$

• As before, we can express w in terms of α_i and make predictions:

$$w = \sum_{i=1}^{m} \alpha_i y^{(i)} x^{(i)}. \qquad w^T x + b = \sum_{i=1}^{m} \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b.$$

- Some differences:
 - Introducing an L_1 regularization turns $0 \le \alpha_i$ to $0 \le \alpha_i \le C$
 - The solution of b^* is modified.
 - The KKT dual complementarity condition: $\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \ge 1$ $\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \le 1$ $0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1.$

Ensemble Methods

Two heads are better than one, or「三個臭皮匠,勝過一個諸葛亮」

Why Ensembles?

- **Statistical:** Multiple minima with same performance (training set too small). Choosing average reduces risk of wrong hypothesis choice.
- **Computational:** get stuck in local minima; results (e.g. decision tree structure + classification) vary strongly depending on training set.
- **Representational:** more expressive than single predictor, e.g.,



Aggregating different linear hypotheses



Linear perceptron hypothesis

Bagging

- BAGG=Bootstrap AGGregation (Breiman 1996)
- Partition your dataset: $\mathscr{L} \to \{\mathscr{L}_1, ..., \mathscr{L}_M\}.$
- Each \mathscr{L}_i is large enough to learn a predictor $g_{\mathscr{L}_i}$.
- For continuous predictors, take average:

$$\hat{g}_{\mathcal{L}}^{A}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^{M} g_{\mathcal{L}_{i}}(\mathbf{x}).$$

• For classification, implement majority rule:

 $\hat{g}_{\mathcal{L}}^{A}(\mathbf{x}) = \arg\max_{j} \sum_{i=1}^{M} I[g_{\mathcal{L}_{i}}(\mathbf{x}) = j], \qquad I[g_{\mathcal{L}_{i}}(x) = j] = \text{indicator function} = 1 \text{ if } g_{\mathcal{L}_{i}}(x) = j$

BAGGing significantly reduces the variance w/o increasing the bias.

Empirical Bootstrapping

-1

-2

-1.0

-0.5

0

x

- Problem of BAGGing: requires a lot of data in each partition.
- Way around: empirical bootstrapping



- Price: increase in bias because we are recycling points.
- Where useful? When faced with unstable learning algorithms (prediction error dominated by variance).

Boosting

• Idea: Give different weights to the predictors.

$$g_A(\mathbf{x}) = \sum_{K=1}^M \alpha_k g_k(\mathbf{x}), \qquad \sum_k \alpha_k = 1.$$

- Initialize equal weights $\alpha_k = 1/M \forall$ points and adjust at each step.
- AdaBoost [Freud, Schapire, 1995]: Consider a classifier for $y \in \{+1, -1\}$
 - Initialize $w_{t=1}(x_n) = 1/N, n = 1, ..., N$.
 - For $t = 1 \cdots$, *T*(desired termination step), do:
 - (1) Select a hypothesis $g_t \in \mathcal{H}$ that minimizes the weighted error

$$\epsilon_t = \sum_{i=1}^N w_t(\mathbf{x}_i) \mathbb{1}(g_t(\mathbf{x}_i) \neq y_i)$$

(2) Let $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$, update the weight for each data \mathbf{x}_n by

$$w_{t+1}(\mathbf{x}_n) \leftarrow w_t(\mathbf{x}_n) \frac{\exp[-\alpha_t y_n g_t(\mathbf{x}_n)]}{Z_t},$$

where $Z_t = \sum_{n=1}^N w_t(\mathbf{x}_n) e^{-\alpha_t y_n g_t(\mathbf{x}_n)}$ ensures all weights add up to unity.

• **Output** $g_A(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^T \alpha_t g_t(\mathbf{x})\right)$

Correctly predicted points have $y_n g_t(x_n) = +1$

 $\epsilon_{\rm t} \sim 0.5$ (as good as guess): do nothing

 $\epsilon_t > 0.5$ (large error): weights enhanced for correctly predicted points

 $\epsilon_t < 0.5$ (small error): enhance weights for wrongly predicted points.

Boosting

- Various variants of boosted algorithms: gradient boosted trees, XGBoost (extreme gradient boosting),
- You are encouraged to explore how these variants work by reading their documentations and experimenting with them:

https://scikit-learn.org/stable/modules/generated/ sklearn.ensemble.GradientBoostingClassifier.html

https://xgboost.readthedocs.io/en/latest/index.html

Random Forest

- Multiple Decision Trees used
- Use different subsets of data for fitting a tree (bagging)
- Use different subsets of features (Outlook, Humidity) to fit the data:
- Trying to tackle high-variance in decision trees
- Reduces correlations between decision trees
- Many variants and applications; recent applications include classifying non-Higgsable gauge groups in F-theory:

https://arxiv.org/pdf/1804.07296.pdf

SUSY Dataset

- Use XGBoost to classify Monte-Carlo simulations into SUSY vs SM. We can compare performance with logistic regression done earlier.
- Feature score to rank significance of features:



https://physics.bu.edu/~pankajm/ML-Notebooks/HTML/NB10_CVIII-XGboost_susy.html

Summary

- Kernel methods
- Soft margins
- Ensemble Methods
 - Bagging
 - Boosting
 - Random Forest