# CSCI567 Machine Learning (Fall 2021)

Prof. Haipeng Luo
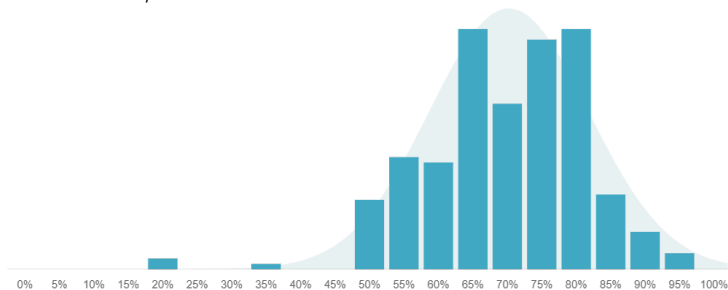
U of Southern California

Oct 21, 2021

# Administration

Quiz 1 grading is done:

- mean: $69.7$, median: $70.5$



- will discuss solutions today

HW3 is due on Tue (Oct 26th)

# Outline

# Outline

# Decision tree

We have seen different ML models for classification/regression:

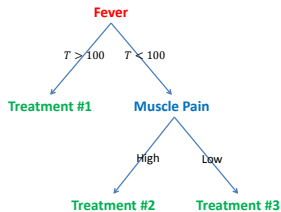- linear models, neural nets and other nonlinear models induced by kernels

**Decision tree** is yet another one:

- nonlinear in general

- works for both classification and regression; we focus on classification

- one key advantage is good interpretability

- used to be very popular; ensemble of trees (i.e. "forest") can still be very effective

- not to be confused with the "tree reduction" in Lec 4

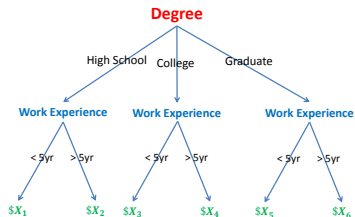# Example
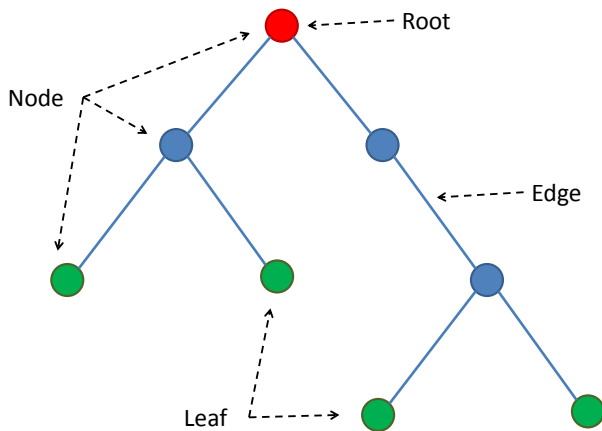
Many decisions are made based on some tree structure

**Medical treatment**



**Salary in a company**

# Tree terminology

# A more abstract example of decision trees

**Input**: $\boldsymbol{x} = (x_1, x_2)$

**Output**: $f(\boldsymbol{x})$ determined naturally by traversing the tree

- start from the root
- test at each node to decide which child to visit next
- finally the leaf gives the prediction $f(\boldsymbol{x})$



For example, $f((\theta_1 - 1, \theta_2 + 1)) = $ B

Complex to formally write down, but easy to represent pictorially or as codes.

# The decision boundary

Corresponds to a classifier with boundaries:

# Parameters

**Parameters** to learn for a decision tree:

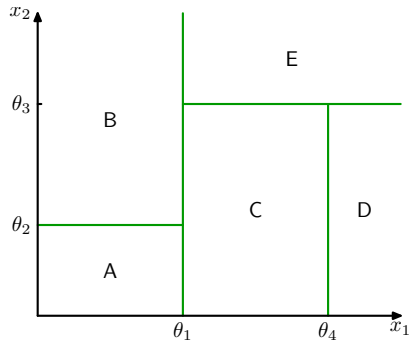- the structure of the tree, such as the depth, #branches, #nodes, etc
  - some of them are sometimes considered as hyperparameters
  - unlike typical neural nets, the structure of a tree is *not fixed in advance, but learned from data*

- the test at each internal node

  - which feature(s) to test on?

  - if the feature is continuous, what threshold $(\theta_1, \theta_2, \ldots)$?



- the value/prediction of the leaves (A, B, ...)

# Learning the parameters

So how do we *learn all these parameters?*

Recall typical approach is to find the parameters that minimize some loss.

This is unfortunately *not feasible for trees*

- For $Z$ nodes, there are roughly $\#\text{features}^Z$ different ways to decide "which feature to test on each node", which is *a lot*.

- enumerating all these configurations to find the one that minimizes some loss is too computationally expensive.

Instead, we turn to some **greedy top-down approach**.
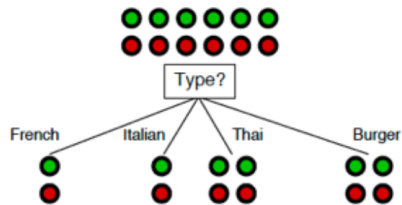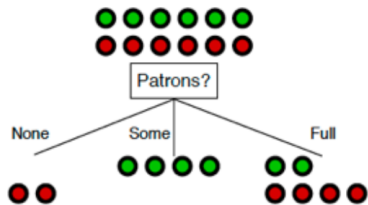
# A running example                         [Russell & Norvig, AIMA]

- predict whether a customer will wait for a table at a restaurant
- 12 training examples
- 10 features (all discrete)

| Example | Attributes | | | | | | | | | | Target |
|---------|------|------|------|------|------|-------|------|------|--------|-------|----------|
|         | $Alt$ | $Bar$ | $Fri$ | $Hun$ | $Pat$ | $Price$ | $Rain$ | $Res$ | $Type$ | $Est$ | $WillWait$ |
| $X_1$   | T | F | F | T | Some | \$\$\$ | F | T | French | 0–10  | T |
| $X_2$   | T | F | F | T | Full | \$     | F | F | Thai   | 30–60 | F |
| $X_3$   | F | T | F | F | Some | \$     | F | F | Burger | 0–10  | T |
| $X_4$   | T | F | T | T | Full | \$     | F | F | Thai   | 10–30 | T |
| $X_5$   | T | F | T | F | Full | \$\$\$ | F | T | French | >60   | F |
| $X_6$   | F | T | F | T | Some | \$\$   | T | T | Italian | 0–10 | T |
| $X_7$   | F | T | F | F | None | \$     | T | F | Burger | 0–10  | F |
| $X_8$   | F | F | F | T | Some | \$\$   | T | T | Thai   | 0–10  | T |
| $X_9$   | F | T | T | F | Full | \$     | T | F | Burger | >60   | F |
| $X_{10}$ | T | T | T | T | Full | \$\$\$ | F | T | Italian | 10–30 | F |
| $X_{11}$ | F | F | F | F | None | \$     | F | F | Thai   | 0–10  | F |
| $X_{12}$ | T | T | T | T | Full | \$     | F | F | Burger | 30–60 | T |

# First step: how to build the root?

I.e., which feature should we test at the root? Examples:



Which split is better?

- intuitively "patrons" is a better feature since it leads to "more pure" or "more certain" children

- how to quantify this intuition?

# Measure of uncertainty of a node

It should be **a function of the distribution of classes**

- e.g. a node with 2 positive and 4 negative examples can be summarized by a distribution $P$ with $P(Y = +1) = 1/3$ and $P(Y = -1) = 2/3$

**Full**

One classic uncertainty measure of a distribution is its *(Shannon) entropy*:

$$H(P) = -\sum_{k=1}^{C} P(Y = k) \log P(Y = k)$$

# Properties of entropy

$$H(P) = -\sum_{k=1}^{\mathsf{C}} P(Y = k) \log P(Y = k)$$

- the base of log can be $2$, $e$ or $10$

- always non-negative

- it's the *smallest codeword length to encode symbols drawn from $P$*

- maximized if $P$ is uniform (max $= \ln \mathsf{C}$): most uncertain case

- minimized if $P$ focuses on one class (min $= 0$): most certain case
    - e.g. $P = (1, 0, \ldots, 0)$
    - $0 \log 0$ is defined naturally as $\lim_{z \to 0+} z \log z = 0$

# Examples of computing entropy

With base $e$ and $4$ classes:

# Another example

Entropy in each child if root tests on "patrons"

For "None" branch

$$-\left( \frac{0}{0+2} \log \frac{0}{0+2} + \frac{2}{0+2} \log \frac{2}{0+2} \right) = 0$$

For "Some" branch

$$-\left( \frac{4}{4+0} \log \frac{4}{4+0} + \frac{0}{4+0} \log \frac{0}{4+0} \right) = 0$$

For "Full" branch

$$-\left( \frac{2}{2+4} \log \frac{2}{2+4} + \frac{4}{2+4} \log \frac{4}{2+4} \right) \approx 0.9$$



*So how good is choosing "patrons" overall?*
Very naturally, we take the weighted average of entropy:

$$\frac{2}{12} \times 0 + \frac{4}{12} \times 0 + \frac{6}{12} \times 0.9 = 0.45$$

# Measure of uncertainty of a split

Suppose we split based on a discrete feature $A$, the uncertainty can be measured by the **conditional entropy**:

$$H(Y \mid A)$$
$$= \sum_a P(A = a) H(Y \mid A = a)$$
$$= \sum_a P(A = a) \left( - \sum_{k=1}^{\mathsf{C}} P(Y \mid A = a) \log P(Y \mid A = a) \right)$$
$$= \sum_a \text{"fraction of example at node } A = a\text{"} \times \text{"entropy at node } A = a\text{"}$$

Pick the feature that leads to the smallest conditional entropy.

# Deciding the root

For "French" branch

$$-\left(\frac{1}{1+1}\log\frac{1}{1+1} + \frac{1}{1+1}\log\frac{1}{1+1}\right) = 1$$

For "Italian" branch

$$-\left(\frac{1}{1+1}\log\frac{1}{1+1} + \frac{1}{1+1}\log\frac{1}{1+1}\right) = 1$$

For "Thai" and "Burger" branches

$$-\left(\frac{2}{2+2}\log\frac{2}{2+2} + \frac{2}{2+2}\log\frac{2}{2+2}\right) = 1$$



The conditional entropy is $\frac{2}{12} \times 1 + \frac{2}{12} \times 1 + \frac{4}{12} \times 1 + \frac{4}{12} \times 1 = 1 > 0.45$

So splitting with "patrons" is better than splitting with "type".

In fact by similar calculation "patrons" is the best split among all features.

We are now done with building the root (this is also called a **stump**).

# Repeat recursively

**Split each child in the same way.**

- but no need to split children "none" and "some": they are pure already and become leaves

- for "full", repeat, focusing on those 6 examples:



| | $Alt$ | $Bar$ | $Fri$ | $Hun$ | $Pat$ | $Price$ | $Rain$ | $Res$ | $Type$ | $Est$ | $WillWait$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $X_1$ | T | F | F | T | Some | \$\$\$ | F | T | French | 0–10 | T |
| $X_2$ | T | F | F | T | Full | \$ | F | F | Thai | 30–60 | F |
| $X_3$ | F | T | F | F | Some | \$ | F | F | Burger | 0–10 | T |
| $X_4$ | T | F | T | T | Full | \$ | F | F | Thai | 10–30 | T |
| $X_5$ | T | F | T | F | Full | \$\$\$ | F | T | French | >60 | F |
| $X_6$ | F | T | F | T | Some | \$\$ | T | T | Italian | 0–10 | T |
| $X_7$ | F | T | F | F | None | \$ | T | F | Burger | 0–10 | F |
| $X_8$ | F | F | F | T | Some | \$\$ | T | T | Thai | 0–10 | T |
| $X_9$ | F | T | T | F | Full | \$ | T | F | Burger | >60 | F |
| $X_{10}$ | T | T | T | T | Full | \$\$\$ | F | T | Italian | 10–30 | F |
| $X_{11}$ | F | F | F | F | None | \$ | F | F | Thai | 0–10 | F |
| $X_{12}$ | T | T | T | T | Full | \$ | F | F | Burger | 30–60 | T |

Again, very easy to interpret.

# Putting it together

**DecisionTreeLearning**(Examples, Features)

- if Examples have the same class, <u>return a leaf with this class</u>

- else if Features is empty, <u>return a leaf with the majority class</u>

- else if Examples is empty, <u>return a leaf with majority class of parent</u>

- else

> find the best feature $A$ to split (e.g. based on conditional entropy)
>
> **Tree** $\leftarrow$ a root with test on $A$
>
> For each value $a$ of $A$:
>
>> **Child** $\leftarrow$ **DecisionTreeLearning**(Examples with $A = a$, Features$\setminus\{A\}$)
>> add **Child** to **Tree** as a new branch

- return **Tree**

# Variants

Popular decision tree algorithms (e.g. C4.5, CART, etc) are all based on this framework.

Variants:

- replace entropy by **Gini impurity**:

$$G(P) = \sum_{k=1}^{C} P(Y = k)(1 - P(Y = k))$$

  meaning: *how often a randomly chosen example would be incorrectly classified if we predict according to another randomly picked example*

- if a feature is continuous, we need to find a threshold that leads to minimum conditional entropy or Gini impurity. *Think about how to do it efficiently.*

# Regularization

If the dataset has no contradiction (i.e. same $x$ but different $y$), the training error of a tree is always zero, which might indicate overfitting.

**Pruning** is a typical way to prevent overfitting for a tree:

- restrict the depth or #nodes

- other more principled approaches

- all make use of a validation set

# Outline

# Introduction

**Boosting**

- is a meta-algorithm, which takes a base algorithm (classification, regression, ranking, etc) as input and boosts its accuracy

- main idea: combine weak "rules of thumb" (e.g. 51% accuracy) to form a highly accurate predictor (e.g. 99% accuracy)

- works very well in practice (especially in combination with trees)

- often is *resistant to overfitting*

- has strong theoretical guarantees

We again focus on binary classification.

# A simple example

**Email spam detection**:

- given a training set like:
    - ("Want to make money fast? ...", **spam**)
    - ("Viterbi Research Gist ...", **not spam**)

- first obtain a classifier by applying a base algorithm, which can be a rather simple/weak one, like decision stumps:
    - e.g. contains the word "money" $\Rightarrow$ spam

- reweight the examples so that "difficult" ones get more attention
    - e.g. spam that doesn't contain the word "money"

- obtain another classifier by applying the same base algorithm:
    - e.g. empty "to address" $\Rightarrow$ spam

- repeat ...

- final classifier is the (weighted) majority vote of all weak classifiers

# The base algorithm

A **base algorithm** $\mathcal{A}$ (also called weak learning algorithm/oracle) takes a training set $S$ weighted by $D$ as input, and outputs classifier $h \leftarrow \mathcal{A}(S, D)$

- this can be any off-the-shelf classification algorithm (e.g. decision trees, logistic regression, neural nets, etc)

- many algorithms can deal with a weighted training set (e.g. for algorithm that minimizes some loss, we can simply replace "total loss" by "weighted total loss")

- even if it's not obvious how to deal with weight directly, we can always resample according to $D$ to create a new unweighted dataset

# Boosting Algorithms

Given:

- a training set $S$

- a base algorithm $\mathcal{A}$

Two things to specify a boosting algorithm:

- how to **reweight** the examples?

- how to **combine** all the weak classifiers?

AdaBoost is one of the most successful boosting algorithms.

# The AdaBoost Algorithm

Given a training set $S$ and a base algorithm $\mathcal{A}$, initialize $D_1$ to be uniform

For $t = 1, \ldots, T$

- obtain a weak classifier $h_t \leftarrow \mathcal{A}(S, D_t)$
- calculate the importance of $h_t$ as

$$\beta_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) \qquad (\beta_t > 0 \Leftrightarrow \epsilon_t < 0.5)$$

  where $\epsilon_t = \sum_{n : h_t(\boldsymbol{x}_n) \neq y_n} D_t(n)$ is the weighted error of $h_t$.

- update distributions

$$D_{t+1}(n) \propto D_t(n) e^{-\beta_t y_n h_t(\boldsymbol{x}_n)} = \begin{cases} D_t(n) e^{-\beta_t} & \text{if } h_t(x_n) = y_n \\ D_t(n) e^{\beta_t} & \text{else} \end{cases}$$

Output the final classifier $H(\boldsymbol{x}) = \text{sgn}\left( \sum_{t=1}^{T} \beta_t h_t(\boldsymbol{x}) \right)$

# Example

10 data points in $\mathbb{R}^2$

The size of + or - indicates the weight, which starts from uniform $D_1$



Base algorithm is decision stump:



$x_1 > \theta_1$

$x_2 > \theta_2$

Observe that *no stump can predict very accurately for this dataset*

# Round 1: $t = 1$



- 3 misclassified (circled): $\epsilon_1 = 0.3 \rightarrow \beta_1 = \frac{1}{2} \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right) \approx 0.42$.

- $D_2$ puts more weights on those examples

# Round 2: $t = 2$



- 3 misclassified (circled): $\epsilon_2 = 0.21 \rightarrow \beta_2 = 0.65$.

- $D_3$ puts more weights on those examples

# Round 3: $t = 3$



- again 3 misclassified (circled): $\epsilon_3 = 0.14 \rightarrow \beta_3 = 0.92$.

# Final classifier: combining 3 classifiers



$$H_{\text{final}} = \text{sign}\left( 0.42 \quad\quad + 0.65 \quad\quad + 0.92 \quad\quad \right)$$

*All data points are now classified correctly*, even though each weak classifier makes 3 mistakes.

# Overfitting

When $T$ is large, the model is very complicated and overfitting can happen



(boosting "stumps" on heart-disease dataset)

# Resistance to overfitting

However, *very often AdaBoost is resistant to overfitting*



(boosting C4.5 on
"letter" dataset)

- test error does not increase, even after 1000 rounds
  - (total size > 2,000,000 nodes)
- test error continues to drop even after training error is zero!

|  | # rounds | | |
|---|---|---|---|
|  | 5 | 100 | 1000 |
| train error | 0.0 | 0.0 | 0.0 |
| test error | 8.4 | 3.3 | 3.1 |

Used to be a mystery, but by now rigorous theory has been developed to explain this phenomenon.

# Why AdaBoost works?

In fact, *AdaBoost also follows the general framework of minimizing some surrogate loss*.

Step 1: the model that AdaBoost considers is

$$\left\{ \operatorname{sgn}\left(f(\cdot)\right) \mid f(\cdot) = \sum_{t=1}^{T} \beta_t h_t(\cdot) \text{ for some } \beta_t \geq 0 \text{ and } h_t \in \mathcal{H} \right\}$$

where $\mathcal{H}$ is the set of models considered by the base algorithm

Step 2: the loss that AdaBoost minimizes is the exponential loss

$$\sum_{n=1}^{N} \exp\left(-y_n f(\boldsymbol{x}_n)\right)$$

# Greedy minimization

Step 3: the way that AdaBoost minimizes exponential loss is by a greedy approach, that is, find $\beta_t, h_t$ one by one for $t = 1, \ldots, T$.

Specifically, let $f_t = \sum_{\tau=1}^{t} \beta_\tau h_\tau$. Suppose we have found $f_{t-1}$, *what should $f_t$ be?*  Greedily, we want to find $\beta_t, h_t$ to minimize

$$\sum_{n=1}^{N} \exp\left(-y_n f_t(\boldsymbol{x}_n)\right) = \sum_{n=1}^{N} \exp\left(-y_n f_{t-1}(\boldsymbol{x}_n)\right) \exp\left(-y_n \beta_t h_t(\boldsymbol{x}_n)\right)$$

$$\propto \sum_{n=1}^{N} D_t(n) \exp\left(-y_n \beta_t h_t(\boldsymbol{x}_n)\right)$$

where the last step is by the definition of weights

$$D_t(n) \propto D_{t-1}(n) \exp\left(-y_n \beta_{t-1} h_{t-1}(\boldsymbol{x}_n)\right) \propto \cdots \propto \exp\left(-y_n f_{t-1}(\boldsymbol{x}_n)\right)$$

# Greedy minimization

So the goal becomes finding $\beta_t, h_t \in \mathcal{H}$ that minimize

$$\sum_{n=1}^{N} D_t(n) \exp\left(-y_n \beta_t h_t(\boldsymbol{x}_n)\right)$$

$$= \sum_{n:y_n \neq h_t(\boldsymbol{x}_n)} D_t(n) e^{\beta_t} + \sum_{n:y_n = h_t(\boldsymbol{x}_n)} D_t(n) e^{-\beta_t}$$

$$= \epsilon_t e^{\beta_t} + (1 - \epsilon_t) e^{-\beta_t} \qquad \left(\text{recall } \epsilon_t = \sum_{n:y_n \neq h_t(\boldsymbol{x}_n)} D_t(n)\right)$$

$$= \epsilon_t (e^{\beta_t} - e^{-\beta_t}) + e^{-\beta_t}$$

It is now clear we should find $h_t$ to minimize its the weighted classification error $\epsilon_t$, *exactly what the base algorithm should do intuitively!*

This greedy step is abstracted out through a base algorithm.

# Greedy minimization

When $h_t$ (and thus $\epsilon_t$) is fixed, we then find $\beta_t$ to minimize

$$\epsilon_t(e^{\beta_t} - e^{-\beta_t}) + e^{-\beta_t}$$

In HW3, you will verify that this exactly gives:

$$\beta_t = \frac{1}{2}\ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$$

Keep doing this greedy minimization gives the AdaBoost algorithm.

# Summary for boosting

Key idea of boosting is to **combine weak predictors into a strong one**.

There are many boosting algorithms; AdaBoost is the most classic one.

AdaBoost is **greedily minimizing the exponential loss**.

AdaBoost is often **resistant to overfitting**.

# Quiz 1 Problem 5 (a)

Consider the following Gaussian/RBF kernel

$$k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(\frac{-\|\boldsymbol{x} - \boldsymbol{x}'\|_2^2}{2}\right). \tag{3}$$

It is known that there exists an infinite-dimensional nonlinear mapping $\phi_{\mathrm{RBF}}$ such that

$$\phi_{\mathrm{RBF}}(\boldsymbol{x})^{\mathrm{T}} \phi_{\mathrm{RBF}}(\boldsymbol{x}') = k(\boldsymbol{x}, \boldsymbol{x}') \tag{4}$$

for any $\boldsymbol{x}$ and $\boldsymbol{x}'$. In this problem, you will investigate a way to approximate this nonlinear mapping $\phi_{\mathrm{RBF}}$.

(a) Consider a nonlinear mapping $\phi_{\boldsymbol{v},b} : \mathbb{R}^D \to \mathbb{R}$ constructed as follows: randomly draw a vector $\boldsymbol{v} \in \mathbb{R}^D$ from the standard Gaussian and a scalar $b$ from the uniform distribution over $[0, \pi]$, then define $\phi_{\boldsymbol{v},b}(\boldsymbol{x}) = \sqrt{2}\cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x} + b)$ for any input feature vector $\boldsymbol{x} \in \mathbb{R}^D$.

For any two feature vectors $\boldsymbol{x}$ and $\boldsymbol{x}'$, prove the following

$$\mathbb{E}\left[\phi_{\boldsymbol{v},b}(\boldsymbol{x})\phi_{\boldsymbol{v},b}(\boldsymbol{x}')\right] = k(\boldsymbol{x}, \boldsymbol{x}') \tag{5}$$

where the expectation is over the randomness of $\boldsymbol{v}$ and $b$, and $k(\cdot, \cdot)$ is defined in Eq. (3). You can directly use the following two identities in your proof:

- trigonometric identity: $2\cos(\alpha)\cos(\beta) = \cos(\alpha - \beta) + \cos(\alpha + \beta)$;
- integral identity: $\mathbb{E}\left[\cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{z})\right] = \exp\left(\frac{-\|\boldsymbol{z}\|_2^2}{2}\right)$ where the expectation is with respect to $\boldsymbol{v}$ randomly drawn from the standard Gaussian. (With this, you do not even need to know what the standard Gaussian is to solve this problem.)

# Quiz 1 Problem 5 (a)

Plugging in the definition of $\phi_{\boldsymbol{v},b}$, we first have

$$\mathbb{E}\left[\phi_{\boldsymbol{v},b}(\boldsymbol{x})\phi_{\boldsymbol{v},b}(\boldsymbol{x}')\right] = 2\mathbb{E}\left[\cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x} + b)\cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x}' + b)\right]. \qquad (1 \text{ point})$$

Using the given trigonometric identity, the above is equal to

$$\mathbb{E}\left[\cos(\boldsymbol{v}^{\mathrm{T}}(\boldsymbol{x} - \boldsymbol{x}')) + \cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x} + \boldsymbol{v}^{\mathrm{T}}\boldsymbol{x}' + 2b)\right]. \qquad (1 \text{ point})$$

For the first term above, directly applying the given integral identity gives

$$\mathbb{E}\left[\cos(\boldsymbol{v}^{\mathrm{T}}(\boldsymbol{x} - \boldsymbol{x}'))\right] = k(\boldsymbol{x}, \boldsymbol{x}'). \qquad (1 \text{ point})$$

For the second term, fixing $\boldsymbol{v}$ and taking the expectation over $b$ shows

$$\begin{aligned}
\mathbb{E}\left[\cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x} + \boldsymbol{v}^{\mathrm{T}}\boldsymbol{x}' + 2b)\right] &= \frac{1}{\pi}\int_0^{\pi}\cos(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x} + \boldsymbol{v}^{\mathrm{T}}\boldsymbol{x}' + 2b)db \\
&= \frac{1}{2\pi}\sin(\boldsymbol{v}^{\mathrm{T}}\boldsymbol{x} + \boldsymbol{v}^{\mathrm{T}}\boldsymbol{x}' + 2b)\Big|_0^{\pi} = 0. \qquad (2 \text{ points})
\end{aligned}$$

This finishes the proof. (The last step can also be argued by symmetry without writing down the integral explicitly.)

# Quiz 1 Problem 5 (b)

(b) Comparing Eq. (4) and Eq. (5), we see that $\phi_{\boldsymbol{v},b}$ can be used as an approximation for $\phi_{\mathrm{RBF}}$. However, using only one sample $(\boldsymbol{v}, b)$ leads to large variance for this approximation. Based on this information, for any given dimension $M > 1$, can you come up with a random nonlinear mapping $\boldsymbol{\phi} : \mathbb{R}^D \to \mathbb{R}^M$, such that it is a better approximation of $\phi_{\mathrm{RBF}}$ satisfying $\mathbb{E}\left[\boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}')\right] = k(\boldsymbol{x}, \boldsymbol{x}')$? Write down your proposal, prove $\mathbb{E}\left[\boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}')\right] = k(\boldsymbol{x}, \boldsymbol{x}')$, and finally explain why it is a better approximation (in one concise sentence). (5 points)

Proposal: $\boldsymbol{\phi}(\boldsymbol{x}) = \left(\frac{1}{\sqrt{M}}\phi_{\boldsymbol{v}_1,b_1}(\boldsymbol{x}), \ldots, \frac{1}{\sqrt{M}}\phi_{\boldsymbol{v}_M,b_M}(\boldsymbol{x})\right)$ where each $(\boldsymbol{v}_j, b_j)$ is an independent sample drawn from the distribution described in the last question.

It satisfies the claimed equality since

$$\mathbb{E}\left[\boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}')\right] = \mathbb{E}\left[\frac{1}{M}\sum_{j=1}^{M}\phi_{\boldsymbol{v}_j,b_j}(\boldsymbol{x})\phi_{\boldsymbol{v}_j,b_j}(\boldsymbol{x}')\right] = \frac{1}{M}\sum_{j=1}^{M}k(\boldsymbol{x}, \boldsymbol{x}') = k(\boldsymbol{x}, \boldsymbol{x}'),$$

where the second step is by Eq. (5). It is a better approximation since using multiple independent samples reduces the variance (by a factor of $1/M$ precisely).

# Quiz 1 Problem 5 (c)

(c) As discussed in Lecture 5, in RBF-kernelized linear regression with training set $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)$, we maintain a weight vector $\boldsymbol{\alpha} = (\boldsymbol{K} + \lambda \boldsymbol{I})^{-1} \boldsymbol{y} \in \mathbb{R}^N$, where $\boldsymbol{K} \in \mathbb{R}^{N \times N}$ is the Gram matrix (such that $K_{n,m} = k(\boldsymbol{x}_n, \boldsymbol{x}_m)$), $\lambda > 0$ is the regularization coefficient, and $\boldsymbol{y} = (y_1, \ldots, y_N)^{\mathrm{T}}$ is the response vector. For a test point $\boldsymbol{x}$, we make a prediction via $\sum_{n=1}^{N} \alpha_n k(\boldsymbol{x}_n, \boldsymbol{x})$. While powerful, this method can be computationally expensive when $N$ is huge.

Based on the nonlinear mapping you proposed in the last question for $M$ much smaller than $N$, describe how you can approximate the kernelized linear regression described above with a much better time and space complexity. You only need to describe what quantities your method maintains, and how it makes a prediction for a test point. (4 points)

The method is simply what we discussed in Lectures 2 and 5: maintain a weight vector $\boldsymbol{w}^* \in \mathbb{R}^M$ as:

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

where the $n$-th row of $\boldsymbol{\Phi} \in \mathbb{R}^{N \times M}$ is $\boldsymbol{\phi}(\boldsymbol{x}_n)^{\mathrm{T}}$. To make a prediction for a test point $\boldsymbol{x}$, simply compute $\boldsymbol{w}^{*\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})$.

Reasoning (NOT required): First, this has better time and space complexity since $M$ is assumed to be much smaller than $N$. Second, based on the discussion in Lecture 5, this is equivalent to kernelized linear regression with Gram matrix $\boldsymbol{\Phi} \boldsymbol{\Phi}^{\mathrm{T}}$, which is a good approximation of $\boldsymbol{K}$ according to the last question.

# Final note

These **random-feature-based methods** are widely successful in practice!

- enjoy the benefit of kernel methods, *without paying for the price*

- Rahimi and Recht won NeurIPS 2017 Test of Time Award for this