CSCI567 Machine Learning (Fall 2025)

Haipeng Luo

University of Southern California

Sep 19, 2025

Administration

Will discuss HW1 solutions in today discussion session.

HW2 will be released next week.

Outline

Review of Last Lecture

Multiclass Classification

3 Kernel methods

Outline

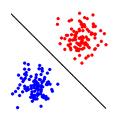
- Review of Last Lecture
- 2 Multiclass Classification
- 3 Kernel methods

Linear classifiers

Linear models for binary classification:

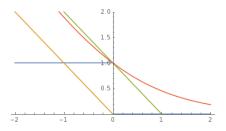
Step 1. Model is the set of separating hyperplanes

$$\mathcal{F} = \{f(\boldsymbol{x}) = \operatorname{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}) \mid \boldsymbol{w} \in \mathbb{R}^{\mathsf{D}}\}$$



Linear classifiers

Step 2. Pick the surrogate loss



- perceptron loss $\ell_{perceptron}(z) = \max\{0, -z\}$ (used in Perceptron)
- hinge loss $\ell_{\text{hinge}}(z) = \max\{0, 1-z\}$ (used in SVM and many others)
- logistic loss $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$ (used in logistic regression)

Linear classifiers

Step 3. Find empirical risk minimizer (ERM):

$$\boldsymbol{w}^* = \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^{\mathsf{D}}} F(\boldsymbol{w}) = \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^{\mathsf{D}}} \frac{1}{N} \sum_{n=1}^{N} \ell(y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n)$$

using

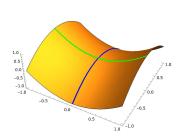
• GD: $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla F(\boldsymbol{w})$

• SGD: $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \tilde{\nabla} F(\boldsymbol{w})$ $(\mathbb{E}[\tilde{\nabla} F(\boldsymbol{w})] = \nabla F(\boldsymbol{w}))$

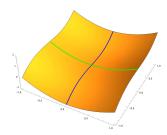
• Newton: $\boldsymbol{w} \leftarrow \boldsymbol{w} - \left(\nabla^2 F(\boldsymbol{w})\right)^{-1} \nabla F(\boldsymbol{w})$

Convergence guarantees of GD/SGD

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or "bad" saddle points (random initialization escapes "good" saddle points)



"good" saddle points



"bad" saddle points

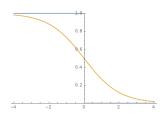
Perceptron and logistic regression

Initialize w = 0 or randomly.

Repeat:

- ullet pick a data point $oldsymbol{x}_n$ uniformly at random (common trick for SGD)
- update parameter:

$$m{w} \leftarrow m{w} + egin{cases} \mathbb{I}[y_n m{w}^{\mathrm{T}} m{x}_n \leq 0] y_n m{x}_n & \text{(Perceptron)} \\ \eta \sigma(-y_n m{w}^{\mathrm{T}} m{x}_n) y_n m{x}_n & \text{(logistic regression)} \end{cases}$$



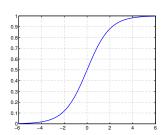
A Probabilistic view of logistic regression

Minimizing logistic loss = MLE for the sigmoid model

$$\boldsymbol{w}^* = \operatorname*{argmin}_{\boldsymbol{w}} \sum_{n=1}^N \ell_{\mathsf{logistic}}(y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n) = \operatorname*{argmax}_{\boldsymbol{w}} \prod_{n=1}^N \mathbb{P}(y_n \mid \boldsymbol{x}_n; \boldsymbol{w})$$

where

$$\mathbb{P}(y \mid \boldsymbol{x}; \boldsymbol{w}) = \sigma(y \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}) = \frac{1}{1 + e^{-y \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}}}$$



Outline

- Review of Last Lecture
- Multiclass Classification
 - Multinomial logistic regression
 - Reduction to binary classification
- Kernel methods

Classification

Recall the setup:

- ullet input (feature vector): $oldsymbol{x} \in \mathbb{R}^{\mathsf{D}}$
- output (label): $y \in [C] = \{1, 2, \dots, C\}$
- ullet goal: learn a mapping $f:\mathbb{R}^{\mathsf{D}} o [\mathsf{C}]$

Examples:

- recognizing digits (C = 10) or letters (C = 26 or 52)
- predicting weather: sunny, cloudy, rainy, etc
- ullet predicting image category: ImageNet dataset (C pprox 20K)

Nearest Neighbor Classifier naturally works for arbitrary C.

Linear models: from binary to multiclass

Step 1: What should a linear model look like for multiclass tasks?

Note: a linear model for binary tasks (switching from $\{-1, +1\}$ to $\{1, 2\}$)

$$f(\boldsymbol{x}) = \begin{cases} 1 & \text{if } \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x} \ge 0 \\ 2 & \text{if } \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x} < 0 \end{cases}$$

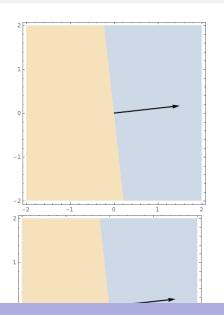
can be written as

$$f(\boldsymbol{x}) = \begin{cases} 1 & \text{if } \boldsymbol{w}_1^{\mathrm{T}} \boldsymbol{x} \geq \boldsymbol{w}_2^{\mathrm{T}} \boldsymbol{x} \\ 2 & \text{if } \boldsymbol{w}_2^{\mathrm{T}} \boldsymbol{x} > \boldsymbol{w}_1^{\mathrm{T}} \boldsymbol{x} \end{cases}$$
$$= \operatorname*{argmax}_{k \in \{1,2\}} \boldsymbol{w}_k^{\mathrm{T}} \boldsymbol{x}$$

for any $\boldsymbol{w}_1, \boldsymbol{w}_2$ s.t. $\boldsymbol{w} = \boldsymbol{w}_1 - \boldsymbol{w}_2$

Think of $w_k^{\mathrm{T}} x$ as a score for class k.

Linear models: from binary to multiclass

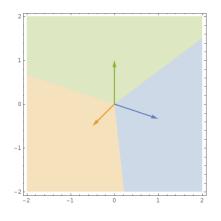


$$w = (\frac{3}{2}, \frac{1}{6}) = w_1 - w_2$$

 $w_1 = (1, -\frac{1}{3})$
 $w_2 = (-\frac{1}{2}, -\frac{1}{2})$

Blue class:

Linear models: from binary to multiclass



$$\mathbf{w}_1 = (1, -\frac{1}{3})$$

 $\mathbf{w}_2 = (-\frac{1}{2}, -\frac{1}{2})$
 $\mathbf{w}_3 = (0, 1)$

Blue class:

$$\{\boldsymbol{x}: 1 = \operatorname{argmax}_k \boldsymbol{w}_k^{\mathrm{T}} \boldsymbol{x}\}$$

• Orange class:

$$\{\boldsymbol{x}: \boldsymbol{2} = \operatorname{argmax}_k \boldsymbol{w}_k^{\mathrm{T}} \boldsymbol{x}\}$$

• Green class:

$$\{ oldsymbol{x} : oldsymbol{3} = \operatorname{argmax}_k oldsymbol{w}_k^{\mathrm{T}} oldsymbol{x} \}$$

Linear models for multiclass classification

$$\mathcal{F} = \left\{ f(oldsymbol{x}) = rgmax_{k \in [\mathsf{C}]} oldsymbol{w}_k^{\mathrm{T}} oldsymbol{x} \mid oldsymbol{w}_1, \dots, oldsymbol{w}_{\mathsf{C}} \in \mathbb{R}^{\mathsf{D}}
ight\}$$

$$= \left\{ f(oldsymbol{x}) = rgmax_{k \in [\mathsf{C}]} oldsymbol{(Woldsymbol{x})_k} \mid oldsymbol{W} \in \mathbb{R}^{\mathsf{C} imes \mathsf{D}}
ight\}$$

Step 2: How do we generalize perceptron/hinge/logistic loss?

This lecture: focus on the more popular logistic loss

Multinomial logistic regression: a probabilistic view

Observe: for binary logistic regression, with $m{w} = m{w}_1 - m{w}_2$:

$$\mathbb{P}(y = 1 \mid \boldsymbol{x}; \boldsymbol{w}) = \sigma(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}}} = \frac{e^{\boldsymbol{w}_{1}^{\mathrm{T}} \boldsymbol{x}}}{e^{\boldsymbol{w}_{1}^{\mathrm{T}} \boldsymbol{x}} + e^{\boldsymbol{w}_{2}^{\mathrm{T}} \boldsymbol{x}}} \propto e^{\boldsymbol{w}_{1}^{\mathrm{T}} \boldsymbol{x}}$$

Naturally, for multiclass:

$$\operatorname{softmax}(\boldsymbol{W}\boldsymbol{x})_k = \mathbb{P}(y = k \mid \boldsymbol{x}; \boldsymbol{W}) = \frac{e^{\boldsymbol{w}_k^{\mathrm{T}}\boldsymbol{x}}}{\sum_{k' \in [\mathsf{C}]} e^{\boldsymbol{w}_{k'}^{\mathrm{T}}\boldsymbol{x}}} \propto e^{\boldsymbol{w}_k^{\mathrm{T}}\boldsymbol{x}}$$

Important operator: softmax function (or really, "softargmax")

For a vector
$$s \in \mathbb{R}^{\mathsf{C}}$$
, $\operatorname{softmax}(s) = \left(\frac{e^{s_1}}{\sum_{k \in [\mathsf{C}]} e^{s_k}}, \cdots, \frac{e^{s_{\mathsf{C}}}}{\sum_{k \in [\mathsf{C}]} e^{s_k}}\right)$

Applying MLE again

Maximize probability of seeing labels y_1, \ldots, y_N given x_1, \ldots, x_N

$$P(\boldsymbol{W}) = \prod_{n=1}^{\mathsf{N}} \mathbb{P}(y_n \mid \boldsymbol{x}_n; \boldsymbol{W}) = \prod_{n=1}^{\mathsf{N}} \frac{e^{\boldsymbol{w}_{y_n}^{\mathsf{T}} \boldsymbol{x}_n}}{\sum_{k \in [\mathsf{C}]} e^{\boldsymbol{w}_k^{\mathsf{T}} \boldsymbol{x}_n}}$$

By taking negative log, this is equivalent to minimizing

$$F(\boldsymbol{W}) = \sum_{n=1}^{N} \ln \left(\frac{\sum_{k \in [C]} e^{\boldsymbol{w}_{k}^{\mathrm{T}} \boldsymbol{x}_{n}}}{e^{\boldsymbol{w}_{y_{n}}^{\mathrm{T}} \boldsymbol{x}_{n}}} \right) = \sum_{n=1}^{N} \ln \left(1 + \sum_{k \neq y_{n}} e^{(\boldsymbol{w}_{k} - \boldsymbol{w}_{y_{n}})^{\mathrm{T}} \boldsymbol{x}_{n}} \right)$$

This is the multiclass logistic loss, a.k.a. cross-entropy loss.

When C = 2, this is the same as binary logistic loss.

Step 3: Optimization

Apply SGD: what is the gradient of

$$F_n(\boldsymbol{W}) = \ln \left(1 + \sum_{k' \neq y_n} e^{(\boldsymbol{w}_{k'} - \boldsymbol{w}_{y_n})^{\mathrm{T}} \boldsymbol{x}_n} \right) ?$$

It's a $C \times D$ matrix. Let's focus on the k-th row:

If $k \neq y_n$:

$$\nabla_{\boldsymbol{w}_{k}^{\mathrm{T}}} F_{n}(\boldsymbol{W}) = \frac{e^{(\boldsymbol{w}_{k} - \boldsymbol{w}_{y_{n}})^{\mathrm{T}} \boldsymbol{x}_{n}}}{1 + \sum_{k' \neq y_{n}} e^{(\boldsymbol{w}_{k'} - \boldsymbol{w}_{y_{n}})^{\mathrm{T}} \boldsymbol{x}_{n}}} \boldsymbol{x}_{n}^{\mathrm{T}} = \mathbb{P}(k \mid \boldsymbol{x}_{n}; \boldsymbol{W}) \boldsymbol{x}_{n}^{\mathrm{T}}$$

else:

$$\nabla_{\boldsymbol{w}_{k}^{\mathrm{T}}} F_{n}(\boldsymbol{W}) = \frac{-\left(\sum_{k' \neq y_{n}} e^{(\boldsymbol{w}_{k'} - \boldsymbol{w}_{y_{n}})^{\mathrm{T}}} \boldsymbol{x}_{n}\right)}{1 + \sum_{k' \neq y_{n}} e^{(\boldsymbol{w}_{k'} - \boldsymbol{w}_{y_{n}})^{\mathrm{T}}} \boldsymbol{x}_{n}} \boldsymbol{x}_{n}^{\mathrm{T}} = (\mathbb{P}(y_{n} \mid \boldsymbol{x}_{n}; \boldsymbol{W}) - 1) \boldsymbol{x}_{n}^{\mathrm{T}}$$

SGD for multinomial logistic regression

Initialize W = 0 (or randomly). Repeat:

- pick $n \in [N]$ uniformly at random
- update the parameters

$$\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \begin{pmatrix} \mathbb{P}(y = 1 \mid \boldsymbol{x}_n; \boldsymbol{W}) \\ \vdots \\ \mathbb{P}(y = y_n \mid \boldsymbol{x}_n; \boldsymbol{W}) - 1 \\ \vdots \\ \mathbb{P}(y = \mathsf{C} \mid \boldsymbol{x}_n; \boldsymbol{W}) \end{pmatrix} \boldsymbol{x}_n^{\mathrm{T}}$$

Think about why the algorithm makes sense intuitively.

A note on prediction

Having learned W, we can either

- ullet make a $extit{deterministic}$ prediction $rgmax_{k \in [\mathsf{C}]} oldsymbol{w}_k^{\mathrm{T}} oldsymbol{x}$
- ullet make a *randomized* prediction drawn from softmax($oldsymbol{W}oldsymbol{x}$)

Generalization of cross-entropy loss

Given a general model class:

$$\mathcal{F} = \left\{ f(oldsymbol{x}) = rgmax_{k \in [\mathsf{C}]} \ s(oldsymbol{x})_k
ight\}$$

where $s: \mathbb{R}^{\mathsf{D}} \to \mathbb{R}^{\mathsf{C}}$ is a "scoring" function.

The *cross-entropy loss* of f for a training sample (x, y) is

$$-\ln\left(\operatorname{softmax}(s(\boldsymbol{x}))_y\right) = -\ln\left(\frac{e^{s(\boldsymbol{x})_y}}{\sum_{k \in [\mathsf{C}]} e^{s(\boldsymbol{x})_k}}\right) = \ln\left(1 + \sum_{k \neq y} e^{s(\boldsymbol{x})_k - s(\boldsymbol{x})_y}\right)$$

Reduce multiclass to binary

Is there an *even more general and simpler approach* to derive multiclass classification algorithms?

Given a binary classification algorithm (any one, not just linear methods), can we turn it to a multiclass algorithm, in a black-box manner?

Yes, there are in fact many ways to do it.

- one-versus-all (one-versus-rest, one-against-all, etc.)
- one-versus-one (all-versus-all, etc.)
- tree-based reduction

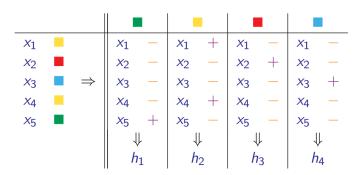
One-versus-all (OvA)

(picture credit: link)

Idea: train C binary classifiers to learn "is class k or not?" for each k.

Training: for each class $k \in [C]$,

- ullet relabel examples with class k as +1, and all others as -1
- ullet train a binary classifier h_k using this new dataset



One-versus-all (OvA)

Prediction: for a new example $oldsymbol{x}$

- ask each h_k : does this belong to class k? (i.e. $h_k(x)$)
- randomly pick among all k's s.t. $h_k(x) = +1$.

Issue: will (probably) make a mistake as long as one of h_k errs.

One-versus-one (OvO)

(picture credit: link)

Idea: train $\binom{\mathsf{C}}{2}$ binary classifiers to learn "is class k or k'?".

Training: for each pair (k, k'),

- ullet relabel examples with class k as +1 and examples with class k' as -1
- discard all other examples
- ullet train a binary classifier $h_{(k,k')}$ using this new dataset

		■ v	s. =	■ v	s. =	■ v	s.	■ v	'S. 📒	■ v	S.	■ v	s. 📙
x_1		<i>x</i> ₁	_					<i>x</i> ₁	_			<i>x</i> ₁	_
<i>x</i> ₂				<i>x</i> ₂	_	<i>x</i> ₂	+					<i>x</i> ₂	+
<i>X</i> 3	\Rightarrow					<i>X</i> 3	_	<i>X</i> 3	+	<i>X</i> 3	_		
<i>X</i> ₄		<i>X</i> ₄	_					<i>X</i> ₄	_			<i>X</i> ₄	_
<i>X</i> 5		<i>X</i> 5	+	<i>X</i> 5	+					<i>X</i> 5	+		
		↓		₩		↓		₩		₩			
		$h_{(1,2)}$		$h_{(1,3)}$		$h_{(3,4)}$		$h_{(4,2)}$		$h_{(1,4)}$		$h_{(3,2)}$	

One-versus-one (OvO)

Prediction: for a new example $oldsymbol{x}$

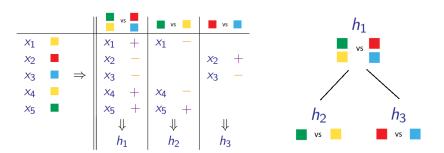
- ask each classifier $h_{(k,k')}$ to vote for either class k or k'
- predict the class with the most votes (break tie in some way)

More robust than one-versus-all, but slower in prediction.

Tree based method

Idea: train \approx C binary classifiers to learn "belongs to which half?".

Training: see pictures



Prediction is also natural, *but is very fast!* (think ImageNet where $C \approx 20K$)

Comparisons

Reduction	training time	prediction time	remark		
OvA	CN	С	not robust		
OvO	(C-1)N	$\mathcal{O}(C^2)$	can achieve very small training error		
Tree	$\mathcal{O}((\log_2 C)N)$	$\mathcal{O}(\log_2C)$	good for "extreme classification"		



Outline

- Review of Last Lecture
- Multiclass Classification
- Kernel methods
 - Motivation
 - Example: Perceptron
 - Kernel Trick
 - Kernelized Perceptron

Motivation

Recall: when linear models are not good enough, we can use a nonlinear feature map $\phi:\mathbb{R}^{\mathsf{D}}\to\mathbb{R}^{\mathsf{M}}$ to transform all x to $\phi(x)$.

Issue: what if M is huge, or even infinity?

Solution: kernel methods

Case study: Perceptron for binary classification

Perceptron

Initialize w=0

Repeat:

- ullet Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}_n) \neq y_n$, update $\boldsymbol{w} \leftarrow \boldsymbol{w} + y_n \boldsymbol{x}_n$

Observation: w is a linear combination of training data

$$\boldsymbol{w} = \sum_{m=1}^{\mathsf{N}} \alpha_m \boldsymbol{x}_m$$

where $\alpha_m = y_m \times \text{number of times } x_m \text{ has been misclassified}$

Dual form of Perceptron

Perceptron (primal form)

Initialize $oldsymbol{w} = oldsymbol{0}$

Repeat:

- ullet Pick a data point index n uniformly at random
- ullet If $\operatorname{sgn}(oldsymbol{w}^{\mathrm{T}}oldsymbol{x}_n)
 eq y_n$, update $oldsymbol{w} \leftarrow oldsymbol{w} + y_noldsymbol{x}_n$

How to update $\alpha_1, \ldots, \alpha_N$ so that $\sum_{m=1}^N \alpha_m x_m \equiv w$ holds always?

Perceptron (dual form)

Initialize $\alpha_m = 0$ for all $m \in [N]$

Repeat:

- ullet Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\sum_{m=1}^{\mathsf{N}} \alpha_m \boldsymbol{x}_{\boldsymbol{m}}^{\mathrm{T}} \boldsymbol{x}_n) \neq y_n$, update $\alpha_n \leftarrow \alpha_n + y_n$

Applying a feature map

Perceptron (primal form with ϕ)

issue: time/space linear in M

Initialize $oldsymbol{w} = oldsymbol{0} \in \mathbb{R}^{oldsymbol{\mathsf{M}}}$

Repeat:

- ullet Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n)) \neq y_n$, update $\boldsymbol{w} \leftarrow \boldsymbol{w} + y_n\boldsymbol{\phi}(\boldsymbol{x}_n)$

Perceptron (dual form with ϕ)

Initialize $\alpha_m = 0$ for all $m \in [N]$

Repeat:

- ullet Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\sum_{m=1}^{N} \alpha_m \phi(\boldsymbol{x}_m)^{\mathrm{T}} \phi(\boldsymbol{x}_n)) \neq y_n$, update $\alpha_n \leftarrow \alpha_n + y_n$

If we can compute $\phi(x_m)^T\phi(x_n)$ without explicitly evaluating $\phi(x_m)$ and $\phi(x_n)$, then time/space is independent of M!

Example

Consider the following polynomial basis $\phi: \mathbb{R}^2 \to \mathbb{R}^3$:

$$\phi(\boldsymbol{x}) = \left(\begin{array}{c} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{array}\right)$$

What is the inner product between $\phi(x)$ and $\phi(x')$?

$$\phi(\mathbf{x})^{\mathrm{T}}\phi(\mathbf{x}') = x_1^2 x_1'^2 + 2x_1 x_2 x_1' x_2' + x_2^2 x_2'^2$$
$$= (x_1 x_1' + x_2 x_2')^2 = (\mathbf{x}^{\mathrm{T}} \mathbf{x}')^2$$

Therefore, the inner product in the new space is simply a function of the inner product in the original space.

Another example

 $\phi: \mathbb{R}^{\mathsf{D}} \to \mathbb{R}^{2\mathsf{D}}$ is parameterized by θ :

$$\phi_{\theta}(\boldsymbol{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_{\mathsf{D}}) \\ \sin(\theta x_{\mathsf{D}}) \end{pmatrix}$$

What is the inner product between $\phi_{ heta}(x)$ and $\phi_{ heta}(x')$?

$$\begin{aligned} \phi_{\theta}(\boldsymbol{x})^{\mathrm{T}} \phi_{\theta}(\boldsymbol{x}') &= \sum_{d=1}^{D} \cos(\theta x_d) \cos(\theta x_d') + \sin(\theta x_d) \sin(\theta x_d') \\ &= \sum_{d=1}^{D} \cos(\theta (x_d - x_d')) \end{aligned} \quad \text{(trigonometric identity)}$$

Once again, the inner product in the new space is a simple function of the features in the original space.

More complicated example

Based on ϕ_{θ} , define $\phi_L : \mathbb{R}^{D} \to \mathbb{R}^{2D(L+1)}$ for some integer L:

$$\phi_L(oldsymbol{x}) = \left(egin{array}{c} \phi_0(oldsymbol{x}) \ \phi_{2rac{2\pi}{L}}(oldsymbol{x}) \ \phi_{2rac{2\pi}{L}}(oldsymbol{x}) \ dots \ \phi_{Lrac{2\pi}{L}}(oldsymbol{x}) \end{array}
ight)$$

What is the inner product between $\phi_L(x)$ and $\phi_L(x')$?

$$egin{aligned} oldsymbol{\phi}_L(oldsymbol{x})^{\mathrm{T}} oldsymbol{\phi}_L(oldsymbol{x}') &= \sum_{\ell=0}^L oldsymbol{\phi}_{rac{2\pi\ell}{L}}(oldsymbol{x})^{\mathrm{T}} oldsymbol{\phi}_{rac{2\pi\ell}{L}}(oldsymbol{x}') \ &= \sum_{\ell=0}^L \sum_{d=1}^{\mathsf{D}} \cos\left(rac{2\pi\ell}{L}(x_d - x_d')
ight) \end{aligned}$$

Infinite dimensional mapping

When $L \to \infty$, even if we cannot compute $\phi(x)$, a vector of *infinite* dimension, we can still compute inner product:

$$\phi_{\infty}(\boldsymbol{x})^{\mathrm{T}}\phi_{\infty}(\boldsymbol{x}') = \int_{0}^{2\pi} \sum_{d=1}^{D} \cos(\theta(x_d - x_d')) d\theta$$
$$= \sum_{d=1}^{D} \frac{\sin(2\pi(x_d - x_d'))}{x_d - x_d'}$$

Again, a simple function of the original features.

Note that using this mapping in linear classification, we are *learning a* weight w with infinite dimension!

Kernel functions

Definition: a function $k : \mathbb{R}^D \times \mathbb{R}^D \to \mathbb{R}$ is called a *kernel function* if there exists a function $\phi : \mathbb{R}^D \to \mathbb{R}^M$ so that for any $x, x' \in \mathbb{R}^D$,

$$k(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}')$$

Can be seen as a kind of similarity measure.

Examples we have seen

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^{\mathrm{T}} \boldsymbol{x}')^{2}$$
$$k(\boldsymbol{x}, \boldsymbol{x}') = \sum_{d=1}^{\mathsf{D}} \frac{\sin(2\pi(x_{d} - x'_{d}))}{x_{d} - x'_{d}}$$

Composing kernels

Creating more kernel functions using the following rules:

If $k_1(\cdot,\cdot)$ and $k_2(\cdot,\cdot)$ are kernels, the followings are kernels too

- conical combination: $\alpha k_1(\cdot,\cdot) + \beta k_2(\cdot,\cdot)$ if $\alpha,\beta \geq 0$
- product: $k_1(\cdot,\cdot)k_2(\cdot,\cdot)$
- exponential: $e^{k(\cdot,\cdot)}$
- o . . .

Verify using the definition of kernel!

Common kernel functions

Two most commonly used kernel functions in practice:

Polynomial kernel

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^{\mathrm{T}} \boldsymbol{x}' + c)^d$$

for $c \ge 0$ and d is a positive integer.

Gaussian kernel or Radial Basis Function (RBF) kernel

$$k(\boldsymbol{x}, \boldsymbol{x}') = e^{-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|_2^2}{2\sigma^2}}$$

for some $\sigma > 0$.

Think about what the corresponding ϕ is for each kernel.

Back to Perceptron

Perceptron (dual form with ϕ)

Initialize $\alpha_m = 0$ for all $m \in [N]$

Repeat:

- ullet Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\sum_{m=1}^{N} \alpha_m \phi(\boldsymbol{x}_m)^{\mathrm{T}} \phi(\boldsymbol{x}_n)) \neq y_n$, update $\alpha_n \leftarrow \alpha_n + y_n$

Instead of choosing $\phi: \mathbb{R}^{D} \to \mathbb{R}^{M}$ explicitly, we choose a kernel function k.

Kernelized Perceptron

Initialize $\alpha_m = 0$ for all $m \in [N]$

Repeat:

- ullet Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\sum_{m=1}^{N} \alpha_m k(\boldsymbol{x}_m, \boldsymbol{x}_n)) \neq y_n$, update $\alpha_n \leftarrow \alpha_n + y_n$

Completely M-independent, becomes a non-parametric method

Gram/kernel matrix

When N is small, can precompute all inner products as a Gram matrix

$$oldsymbol{K} = \left(egin{array}{cccc} k(oldsymbol{x}_1, oldsymbol{x}_1) & k(oldsymbol{x}_1, oldsymbol{x}_2) & \cdots & k(oldsymbol{x}_1, oldsymbol{x}_N) \ k(oldsymbol{x}_2, oldsymbol{x}_1) & k(oldsymbol{x}_2, oldsymbol{x}_2) & \cdots & k(oldsymbol{x}_2, oldsymbol{x}_N) \ k(oldsymbol{x}_N, oldsymbol{x}_1) & k(oldsymbol{x}_N, oldsymbol{x}_2) & \cdots & k(oldsymbol{x}_N, oldsymbol{x}_N) \end{array}
ight) = oldsymbol{\Phi}^{\mathrm{T}}$$

Recall:
$$oldsymbol{\Phi} = \left(egin{array}{c} oldsymbol{\phi}(oldsymbol{x}_1)^{\mathrm{T}} \ oldsymbol{\phi}(oldsymbol{x}_N)^{\mathrm{T}} \end{array}
ight) \in \mathbb{R}^{N imes M}$$

Gram matrix vs covariance matrix

	dimensions	entry (i,j)	property
$\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}$	$N \times N$	$oldsymbol{\phi}(oldsymbol{x}_i)^{ ext{T}}oldsymbol{\phi}(oldsymbol{x}_j)$	both are symmetric and
$\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}$	$M \times M$	$\sum_{n=1}^{N} \phi(\boldsymbol{x}_n)_i \phi(\boldsymbol{x}_n)_j$	positive semidefinite

Mercer Theorem

 $k: \mathbb{R}^{D} \times \mathbb{R}^{D} \to \mathbb{R}$ is a kernel function if and only if the Gram matrix K for any N and any $x_1, x_2, ..., x_N$ is positive semidefinite.

useful for showing that a function is not a kernel

Example:

$$k(x, x') = ||x - x'||_2^2$$

is not a kernel, why?

If it is a kernel, the kernel matrix for two data points x_1 and x_2 :

$$K = \begin{pmatrix} 0 & \|x_1 - x_2\|_2^2 \\ \|x_1 - x_2\|_2^2 & 0 \end{pmatrix}$$

must be positive semidefinite, but it is not (contradiction).

Kernelizing ML algorithms

Many other ML algorithms can be kernelized:

- nearest neighbor classifier
- linear regression
- logistic regression
- SVM
- o . . .

Key idea: rewrite the algorithm so that its dependence on the transformed dataset Φ is only through the Gram matrix $K=\Phi\Phi^{\mathrm{T}}$.