Kernel Methods: From Basics to Modern Applications

Danica J. Sutherland (she/her) Computer Science, University of British Columbia Data Science Summer School, January 2021

• Machine learning!

• Machine learning! ...but how do we actually do it?

- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$

- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$

- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$



- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$
- Extend *x*...

$$f(x) = w^{\mathsf{T}}(1,x,x^2) = w^{\mathsf{T}}\phi(x)$$

- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$
- Extend *x*...

$$f(x) = w^{\mathsf{T}}(1,x,x^2) = w^{\mathsf{T}}\phi(x)$$



- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$
- Extend *x*...

$$f(x) = w^{\mathsf{T}}(1,x,x^2) = w^{\mathsf{T}}\phi(x)$$



- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$
- Extend *x*...

$$f(x) = w^{\mathsf{T}}(1,x,x^2) = w^{\mathsf{T}}\phi(x)$$

- Kernels are basically a way to study doing this with any, potentially very complicated, ϕ

- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$
- Extend *x*...

$$f(x) = w^{\mathsf{T}}(1,x,x^2) = w^{\mathsf{T}}\phi(x)$$

- Kernels are basically a way to study doing this with any, potentially very complicated, ϕ
- Convenient way to make models on documents, graphs, videos, datasets, ...

- Machine learning! ...but how do we actually do it?
- Linear models! $f(x) = w_0 + wx$, $\hat{y}(x) = \mathrm{sign}(f(x))$
- Extend *x*...

$$f(x) = w^{\mathsf{T}}(1,x,x^2) = w^{\mathsf{T}}\phi(x)$$

- Kernels are basically a way to study doing this with any, potentially very complicated, ϕ
- Convenient way to make models on documents, graphs, videos, datasets, ...
- ϕ will live in a *reproducing kernel Hilbert space*

• A complete (real or complex) inner product space.

• A complete (real or complex) inner product space.

- A complete (real or complex) inner product space.
- Inner product space: a vector space with an **inner product**:
 - $\langle lpha_1 f_1 + lpha_2 f_2, g
 angle_{\mathcal{H}} = lpha_1 \langle f_1, g
 angle_{\mathcal{H}} + lpha_2 \langle f_2, g
 angle_{\mathcal{H}}$

•
$$\langle f,g
angle_{\mathcal{H}}=\langle g,f
angle_{\mathcal{H}}$$

• $\langle f,f
angle_{\mathcal{H}}>0$ for f
eq 0, $\langle 0,0
angle_{\mathcal{H}}=0$

- A complete (real or complex) inner product space.
- Inner product space: a vector space with an **inner product**:
 - $\langle lpha_1 f_1 + lpha_2 f_2, g
 angle_{\mathcal{H}} = lpha_1 \langle f_1, g
 angle_{\mathcal{H}} + lpha_2 \langle f_2, g
 angle_{\mathcal{H}}$

•
$$\langle f,g
angle_{\mathcal{H}}=\langle g,f
angle_{\mathcal{H}}$$

• $\langle f,f
angle_{\mathcal{H}}>0$ for f
eq 0, $\langle 0,0
angle_{\mathcal{H}}=0$

Induces a $\operatorname{norm}: \|f\|_{\mathcal{H}} = \sqrt{\langle f, f
angle_{\mathcal{H}}}$

- A complete (real or complex) inner product space.
- Inner product space: a vector space with an **inner product**:
 - $\bullet \ \langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle_{\mathcal{H}} = \alpha_1 \langle f_1, g \rangle_{\mathcal{H}} + \alpha_2 \langle f_2, g \rangle_{\mathcal{H}}$

•
$$\langle f,g
angle_{\mathcal{H}}=\langle g,f
angle_{\mathcal{H}}$$

• $\langle f,f
angle_{\mathcal{H}}>0$ for f
eq 0, $\langle 0,0
angle_{\mathcal{H}}=0$

Induces a $\operatorname{norm}: \|f\|_{\mathcal{H}} = \sqrt{\langle f, f
angle_{\mathcal{H}}}$

• Complete: "well-behaved" (Cauchy sequences have limits in \mathcal{H})

- Call our domain ${\mathcal X}$, some set
 - \mathbb{R}^d , functions, distributions of graphs of images, ...

- Call our domain ${\mathcal X}$, some set
 - \mathbb{R}^d , functions, distributions of graphs of images, ...
- $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a kernel on \mathcal{X} if there exists a Hilbert space \mathcal{H} and a *feature map* $\phi: \mathcal{X} \to \mathcal{H}$ so that

$$k(x,y) = \langle \phi(x), \phi(y)
angle_{\mathcal{H}}$$

- Call our domain ${\mathcal X}$, some set
 - \mathbb{R}^d , functions, distributions of graphs of images, ...
- $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a kernel on \mathcal{X} if there exists a Hilbert space \mathcal{H} and a *feature map* $\phi: \mathcal{X} \to \mathcal{H}$ so that

$$k(x,y)=\langle \phi(x),\phi(y)
angle_{\mathcal{H}}$$

• Roughly, $m{k}$ is a notion of "similarity" between inputs

- Call our domain ${\mathcal X}$, some set
 - \mathbb{R}^d , functions, distributions of graphs of images, ...
- $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a kernel on \mathcal{X} if there exists a Hilbert space \mathcal{H} and a *feature map* $\phi: \mathcal{X} \to \mathcal{H}$ so that

$$k(x,y) = \langle \phi(x), \phi(y)
angle_{\mathcal{H}}$$

- Roughly, $m{k}$ is a notion of "similarity" between inputs
- Linear kernel on \mathbb{R}^d : $k(x,y) = \langle x,y
 angle_{\mathbb{R}^d}$

• Scaling: if $\gamma \geq 0$, $k_\gamma(x,y) = \gamma k(x,y)$ is a kernel

• Scaling: if $\gamma \ge 0$, $k_{\gamma}(x,y) = \gamma k(x,y)$ is a kernel • $k_{\gamma}(x,y) = \gamma \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \sqrt{\gamma} \phi(x), \sqrt{\gamma} \phi(y) \rangle_{\mathcal{H}}$

- Scaling: if $\gamma \ge 0$, $k_{\gamma}(x, y) = \gamma k(x, y)$ is a kernel • $k_{\gamma}(x, y) = \gamma \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \sqrt{\gamma} \phi(x), \sqrt{\gamma} \phi(y) \rangle_{\mathcal{H}}$
- Sum: $k_+(x,y)=k_1(x,y)+k_2(x,y)$ is a kernel

- Scaling: if $\gamma \ge 0$, $k_{\gamma}(x, y) = \gamma k(x, y)$ is a kernel • $k_{\gamma}(x, y) = \gamma \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \sqrt{\gamma} \phi(x), \sqrt{\gamma} \phi(y) \rangle_{\mathcal{H}}$
- Sum: $k_+(x,y) = k_1(x,y) + k_2(x,y)$ is a kernel • $k_+(x,y) = \left\langle \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix}, \begin{bmatrix} \phi_1(y) \\ \phi_2(y) \end{bmatrix} \right\rangle_{\mathcal{H}_1 \oplus \mathcal{H}_2}$

- Scaling: if $\gamma \ge 0$, $k_{\gamma}(x,y) = \gamma k(x,y)$ is a kernel • $k_{\gamma}(x,y) = \gamma \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \sqrt{\gamma} \phi(x), \sqrt{\gamma} \phi(y) \rangle_{\mathcal{H}}$
- Sum: $k_+(x,y) = k_1(x,y) + k_2(x,y)$ is a kernel • $k_+(x,y) = \left\langle \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix}, \begin{bmatrix} \phi_1(y) \\ \phi_2(y) \end{bmatrix} \right\rangle_{\mathcal{H}_1 \oplus \mathcal{H}_2}$
- Is $k_1(x,y)-k_2(x,y)$ necessarily a kernel?

- Scaling: if $\gamma \ge 0$, $k_{\gamma}(x,y) = \gamma k(x,y)$ is a kernel • $k_{\gamma}(x,y) = \gamma \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} = \langle \sqrt{\gamma} \phi(x), \sqrt{\gamma} \phi(y) \rangle_{\mathcal{H}}$
- Sum: $k_+(x,y) = k_1(x,y) + k_2(x,y)$ is a kernel • $k_+(x,y) = \left\langle \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix}, \begin{bmatrix} \phi_1(y) \\ \phi_2(y) \end{bmatrix} \right\rangle_{\mathcal{H}_1 \oplus \mathcal{H}_2}$
- Is $k_1(x,y) k_2(x,y)$ necessarily a kernel? • Take $k_1(x,y) = 0$, $k_2(x,y) = xy$, $x \neq 0$.
 - Then $k_1(x,x)-k_2(x,x)=-x^2<0$
 - But $k(x,x) = \|\phi(x)\|_{\mathcal{H}}^2 \geq 0.$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

• Equivalently: *kernel matrix* $oldsymbol{K}$ is PSD

$$K := egin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \ dots & dots & \ddots & dots \ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix}$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

$$\sum_{i=1}^n \sum_{j=1}^n \langle a_i \phi(x_i), a_j \phi(x_j)
angle_{\mathcal{H}}$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

$$\sum_{i=1}^n\sum_{j=1}^n\langle a_i\phi(x_i),a_j\phi(x_j)
angle_{\mathcal{H}}=\left\langle\sum_{i=1}^na_i\phi(x_i),\sum_{j=1}^na_j\phi(x_j)
ight
angle_{\mathcal{H}}$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

$$egin{aligned} &\sum_{i=1}^n\sum_{j=1}^n\langle a_i\phi(x_i),a_j\phi(x_j)
angle_{\mathcal{H}}=\left\langle\sum_{i=1}^na_i\phi(x_i),\sum_{j=1}^na_j\phi(x_j)
ight
angle_{\mathcal{H}}\ &=\left\|\sum_{i=1}^na_i\phi(x_i)
ight\|_{\mathcal{H}}^2 \end{aligned}$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

$$egin{aligned} &\sum_{i=1}^n\sum_{j=1}^n\langle a_i\phi(x_i),a_j\phi(x_j)
angle_{\mathcal{H}}=\left\langle\sum_{i=1}^na_i\phi(x_i),\sum_{j=1}^na_j\phi(x_j)
ight
angle_{\mathcal{H}}\ &=\left\|\sum_{i=1}^na_i\phi(x_i)
ight\|_{\mathcal{H}}^2\geq 0 \end{aligned}$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

• A symmetric function $k: \mathcal{X} imes \mathcal{X} o \mathbb{R}$ is *positive semi-definite* (*psd*) if for all $n \geq 1, a_1, \ldots, a_n \in \mathbb{R}^n$, $x_1, \ldots, x_n \in \mathcal{X}^n$,

$$\sum_{i=1}^n\sum_{j=1}^na_ia_jk(x_i,x_j)\geq 0$$

- Hilbert space kernels are psd
- psd functions are Hilbert space kernels
 - Moore-Aronszajn Theorem; we'll come back to this
- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd

• Limits: if $k_\infty(x,y) = \lim_{n o\infty} k_n(x,y)$ exists, k_∞ is psd • $\lim_{n o\infty} \sum_{i=1}^n \sum_{j=1}^n a_i a_j k_n(x_i,x_j) \ge 0$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
 - Let $V \sim \mathcal{N}(0,K_1)$, $W \sim \mathcal{N}(0,K_2)$ be independent
 - $\operatorname{Cov}(V_iW_i,V_jW_j)=\operatorname{Cov}(V_i,V_j)\operatorname{Cov}(W_i,W_j)=k_{ imes}(x_i,x_j)$
 - Covariance matrices are psd, so $k_{ imes}$ is too

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y) = k(x,y)^n$ is pd for any integer $n \geq 0$ $x^{\mathsf{T}}y$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y) = k(x,y)^n$ is pd for any integer $n \geq 0$ $x^{\mathsf{T}}y + c$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y) = k(x,y)^n$ is pd for any integer $n \geq 0$ $ig(x^{\mathsf{T}}y+cig)^n$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y) = k(x,y)^n$ is pd for any integer $n \ge 0$

 $(x^{ op}y+c)^n$, the polynomial kernel

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd

•
$$k_{\exp}(x,y) = \lim_{N o \infty} \sum_{n=0}^N rac{1}{n!} k(x,y)^n$$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

• Use the feature map $x\mapsto f(x)\phi(x)$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

$$x^{\mathsf{T}}y$$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

$$rac{1}{\sigma^2}x^{\mathsf{T}}y$$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

$$\exp\left(\frac{1}{\sigma^2} x^\mathsf{T} y\right)$$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

$$\exp\Big(-rac{1}{2\sigma^2}\|x\|^2\Big)\exp\Big(rac{1}{\sigma^2}x^{\mathsf{ extsf{ imes}}}y\Big)\exp\Big(-rac{1}{2\sigma^2}\|y\|^2\Big)$$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

$$\exp \Big(-rac{1}{2\sigma^2} \|x\|^2 \Big) \exp \Big(rac{1}{\sigma^2} x^{\mathsf{T}} y \Big) \exp \Big(-rac{1}{2\sigma^2} \|y\|^2 \Big)$$

$$= \exp \Big(- rac{1}{2\sigma^2} ig[\|x\|^2 - 2x^{\mathsf{T}}y + \|y\|^2 ig] \Big)$$

- Limits: if $k_\infty(x,y) = \lim_{n o \infty} k_n(x,y)$ exists, k_∞ is psd
- Products: $k_ imes(x,y)=k_1(x,y)k_2(x,y)$ is psd
- Powers: $k_n(x,y)=k(x,y)^n$ is pd for any integer $n\geq 0$
- Exponents: $k_{ ext{exp}}(x,y) = \exp(k(x,y))$ is pd
- If $f:\mathcal{X}
 ightarrow\mathbb{R}$, $k_f(x,y)=f(x)k(x,y)f(y)$ is pd

$$egin{aligned} &\exp\left(-rac{1}{2\sigma^2}\|x\|^2
ight)\exp\left(rac{1}{\sigma^2}x^{\mathsf{T}}y
ight)\exp\left(-rac{1}{2\sigma^2}\|y\|^2
ight)\ &=\exp\left(-rac{\|x-y\|^2}{2\sigma^2}
ight)$$
, the Gaussian kernel

$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$

$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$



• Recall original motivating example with

$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$

• Kernel is $k(x,y) = \langle \phi(x), \phi(y)
angle_{\mathcal{H}} = 1 + xy + x^2y^2$



$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$

- Kernel is $k(x,y) = \langle \phi(x), \phi(y)
 angle_{\mathcal{H}} = 1 + xy + x^2y^2$
- Classifier based on linear $f(x) = \langle w, \phi(x)
 angle_{\mathcal{H}}$



$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$

- Kernel is $k(x,y) = \langle \phi(x), \phi(y)
 angle_{\mathcal{H}} = 1 + xy + x^2y^2$
- Classifier based on linear $f(x) = \langle w, \phi(x)
 angle_{\mathcal{H}}$
- $f(\cdot)$ is the function f itself, represented by a vector in \mathbb{R}^3 $f(x) \in \mathbb{R}$ is the function evaluated at a point x

$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$

- Kernel is $k(x,y) = \langle \phi(x), \phi(y)
 angle_{\mathcal{H}} = 1 + xy + x^2y^2$
- Classifier based on linear $f(x) = \langle w, \phi(x)
 angle_{\mathcal{H}}$
- $f(\cdot)$ is the function f itself, represented by a vector in \mathbb{R}^3 $f(x) \in \mathbb{R}$ is the function evaluated at a point x
- Elements of $\mathcal H$ correspond to functions, $f:\mathcal X o\mathbb R$

$$\mathcal{X}=\mathbb{R} \qquad \phi(x)=(1,x,x^2)\in \mathbb{R}^3$$

- Kernel is $k(x,y) = \langle \phi(x), \phi(y)
 angle_{\mathcal{H}} = 1 + xy + x^2y^2$
- Classifier based on linear $f(x) = \langle w, \phi(x)
 angle_{\mathcal{H}}$
- $f(\cdot)$ is the function f itself, represented by a vector in \mathbb{R}^3 $f(x) \in \mathbb{R}$ is the function evaluated at a point x
- Elements of $\mathcal H$ correspond to functions, $f:\mathcal X o\mathbb R$
- Reproducing prop.: $f(x) = \langle f(\cdot), \phi(x)
 angle_{\mathcal{H}}$ for $f \in \mathcal{H}$

Reproducing kernel Hilbert space (RKHS)

• Every psd kernel k on \mathcal{X} defines a (unique) Hilbert space, its RKHS \mathcal{H} , and a map $\phi:\mathcal{X}\to\mathcal{H}$ where

$$k(x,y)=\langle \phi(x),\phi(y)
angle_{\mathcal{H}}$$

• Elements $f \in \mathcal{H}$ are functions on \mathcal{X} , with

$$f(x) = \langle f, \phi(x)
angle_{\mathcal{H}}$$

- Combining the two, we sometimes write $k(x,\cdot)=\phi(x)$

Reproducing kernel Hilbert space (RKHS)

• Every psd kernel k on \mathcal{X} defines a (unique) Hilbert space, its RKHS \mathcal{H} , and a map $\phi:\mathcal{X}\to\mathcal{H}$ where

$$k(x,y) = \langle \phi(x), \phi(y)
angle_{\mathcal{H}}$$

• Elements $f \in \mathcal{H}$ are functions on \mathcal{X} , with

$$f(x) = \langle f, \phi(x)
angle_{\mathcal{H}}$$

- Combining the two, we sometimes write $k(x,\cdot)=\phi(x)$
- $k(x, \cdot)$ is the **evaluation functional** An RKHS is defined by it being *continuous*, or

$$|f(x)| \leq M_x \|f\|_{\mathcal{H}}$$

- Building \mathcal{H} for a given psd k:
 - Start with $\mathcal{H}_0 = \mathrm{span}(\{k(x, \cdot): x \in \mathcal{X}\})$

- Building \mathcal{H} for a given psd k:
 - Start with $\mathcal{H}_0 = \mathrm{span}(\{k(x, \cdot): x \in \mathcal{X}\})$
 - Define $\langle \cdot, \cdot
 angle_{\mathcal{H}_0}$ from $\langle k(x, \cdot), k(y, \cdot)
 angle_{\mathcal{H}_0} = k(x, y)$

- Building \mathcal{H} for a given psd k:
 - Start with $\mathcal{H}_0 = \mathrm{span}(\{k(x, \cdot): x \in \mathcal{X}\})$
 - Define $\langle \cdot, \cdot
 angle_{\mathcal{H}_0}$ from $\langle k(x, \cdot), k(y, \cdot)
 angle_{\mathcal{H}_0} = k(x, y)$
 - Take $\mathcal H$ to be completion of $\mathcal H_0$ in the metric from $\langle\cdot,\cdot
 angle_{\mathcal H_0}$

- Building \mathcal{H} for a given psd k:
 - Start with $\mathcal{H}_0 = \mathrm{span}(\{k(x, \cdot): x \in \mathcal{X}\})$
 - Define $\langle \cdot, \cdot
 angle_{\mathcal{H}_0}$ from $\langle k(x, \cdot), k(y, \cdot)
 angle_{\mathcal{H}_0} = k(x, y)$
 - Take $\mathcal H$ to be completion of $\mathcal H_0$ in the metric from $\langle\cdot,\cdot
 angle_{\mathcal H_0}$
 - Get that the reproducing property holds for $k(x,\cdot)$ in ${\mathcal H}$
Moore-Aronszajn Theorem

- Building \mathcal{H} for a given psd k:
 - Start with $\mathcal{H}_0 = \mathrm{span}(\{k(x, \cdot): x \in \mathcal{X}\})$
 - Define $\langle \cdot, \cdot
 angle_{\mathcal{H}_0}$ from $\langle k(x, \cdot), k(y, \cdot)
 angle_{\mathcal{H}_0} = k(x, y)$
 - Take $\mathcal H$ to be completion of $\mathcal H_0$ in the metric from $\langle\cdot,\cdot
 angle_{\mathcal H_0}$
 - Get that the reproducing property holds for $k(x,\cdot)$ in ${\mathcal H}$
 - Can also show uniqueness

Moore-Aronszajn Theorem

- Building \mathcal{H} for a given psd k:
 - Start with $\mathcal{H}_0 = \mathrm{span}(\{k(x, \cdot): x \in \mathcal{X}\})$
 - Define $\langle \cdot, \cdot
 angle_{\mathcal{H}_0}$ from $\langle k(x, \cdot), k(y, \cdot)
 angle_{\mathcal{H}_0} = k(x, y)$
 - Take $\mathcal H$ to be completion of $\mathcal H_0$ in the metric from $\langle\cdot,\cdot
 angle_{\mathcal H_0}$
 - Get that the reproducing property holds for $k(x,\cdot)$ in ${\mathcal H}$
 - Can also show uniqueness
- Theorem: $m{k}$ is psd iff it's the reproducing kernel of an RKHS

A quick check: linear kernels

•
$$k(x,y) = x^{\mathsf{T}} y$$
 on $\mathcal{X} = \mathbb{R}^d$

• If
$$f(y) = \sum_{i=1}^n a_i k(x_i,y)$$
, then $f(y) = [\sum_{i=1}^\infty a_i x_i]^{\mathsf{T}} y$

- Closure doesn't add anything here, since \mathbb{R}^d is closed
- So, linear kernel gives you RKHS of linear functions

$$ullet \ \|f\|_{\mathcal{H}} = \sqrt{\sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j)} = \|\sum_{i=1}^n a_i x_i\|_{\mathcal{H}}$$

$$k(x,y) = \exp(rac{1}{2\sigma^2} \|x-y\|^2)$$



$$k(x,y) = \exp(rac{1}{2\sigma^2} \|x-y\|^2)$$

$$k(x,y)=\exp(rac{1}{2\sigma^2}\|x-y\|^2)$$



$$k(x,y) = \exp(rac{1}{2\sigma^2} \|x-y\|^2)$$



$$k(x,y) = \exp(rac{1}{2\sigma^2} \|x-y\|^2)$$



$$k(x,y)=\exp(rac{1}{2\sigma^2}\|x-y\|^2)$$

- ${\cal H}$ is infinite-dimensional
- Functions in $\mathcal H$ are bounded: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal H} \leq \sqrt{k(x, x)} \| f \|_{\mathcal H} = \| f \|_{\mathcal H}$



$$k(x,y)=\exp(rac{1}{2\sigma^2}\|x-y\|^2)$$

- ${\cal H}$ is infinite-dimensional
- Functions in $\mathcal H$ are bounded: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal H} \leq \sqrt{k(x, x)} \| f \|_{\mathcal H} = \| f \|_{\mathcal H}$
- Choice of σ controls how fast functions can vary:

$$f(x+t)-f(x)\leq \|k(x+t,\cdot)-k(x',\cdot)\|_{\mathcal{H}}\|f\|_{\mathcal{H}} \ \|k(x+t,\cdot)-k(x,\cdot)\|_{\mathcal{H}}^2 = 2-2k(x,x+t) = 2-2\expigg(-rac{\|t\|^2}{2\sigma^2}igg)$$



$$k(x,y) = \exp(rac{1}{2\sigma^2} \|x-y\|^2)$$

- ${\cal H}$ is infinite-dimensional
- Functions in $\mathcal H$ are bounded: $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal H} \leq \sqrt{k(x, x)} \| f \|_{\mathcal H} = \| f \|_{\mathcal H}$
- Choice of σ controls how fast functions can vary:

$$f(x+t)-f(x)\leq \|k(x+t,\cdot)-k(x',\cdot)\|_{\mathcal{H}}\|f\|_{\mathcal{H}} \ \|k(x+t,\cdot)-k(x,\cdot)\|_{\mathcal{H}}^2\|f\|_{\mathcal{H}} = 2-2k(x,x+t) = 2-2\expigg(-rac{\|t\|^2}{2\sigma^2}igg)$$

• Can say lots more with Fourier properties

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

Linear kernel gives normal ridge regression:

$$\hat{f}\left(x
ight) = \hat{w}^{\mathsf{T}}x; \hspace{1em} \hat{w} = rgmin_{w\in \mathbb{R}^d} \sum_{i=1}^n (w^{\mathsf{T}}x_i - y_i)^2 + \lambda \|w\|^2$$

Nonlinear kernels will give nonlinear regression!

$$\hat{f} = rgmin_{f\in\mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

How to find \hat{f} ?

$$\hat{f} = rgmin_{f\in\mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

How to find \hat{f} ? Representer Theorem

• Let $\mathcal{H}_X = ext{span}\{k(x_i,\cdot)\}_{i=1}^n$ \mathcal{H}_\perp its orthogonal complement in \mathcal{H}

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

- Let $\mathcal{H}_X = ext{span}\{k(x_i,\cdot)\}_{i=1}^n$ \mathcal{H}_\perp its orthogonal complement in \mathcal{H}
- Decompose $f=f_X+f_\perp$ with $f_\mathcal{X}\in\mathcal{H}_X$, $f_\perp\in\mathcal{H}_\perp$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

- Let $\mathcal{H}_X = ext{span}\{k(x_i,\cdot)\}_{i=1}^n$ \mathcal{H}_\perp its orthogonal complement in \mathcal{H}
- Decompose $f=f_X+f_\perp$ with $f_\mathcal{X}\in\mathcal{H}_X$, $f_\perp\in\mathcal{H}_\perp$
- $\bullet \ f(x_i) = \langle f_X + f_{\bot}, k(x_i, \cdot) \rangle_{\mathcal{H}} = \langle f_X, k(x_i, \cdot) \rangle_{\mathcal{H}}$

$$\hat{f} = rgmin_{f\in\mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

- Let $\mathcal{H}_X = ext{span}\{k(x_i,\cdot)\}_{i=1}^n$ \mathcal{H}_\perp its orthogonal complement in \mathcal{H}
- Decompose $f=f_X+f_\perp$ with $f_\mathcal{X}\in\mathcal{H}_X$, $f_\perp\in\mathcal{H}_\perp$
- $\bullet \ f(x_i) = \langle f_X + f_{\bot}, k(x_i, \cdot) \rangle_{\mathcal{H}} = \langle f_X, k(x_i, \cdot) \rangle_{\mathcal{H}}$
- $\| \| \|_{\mathcal{H}}^2 = \| f_X \|_{\mathcal{H}}^2 + \| f_\perp \|_{\mathcal{H}}^2$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

- Let $\mathcal{H}_X = ext{span}\{k(x_i,\cdot)\}_{i=1}^n$ \mathcal{H}_\perp its orthogonal complement in \mathcal{H}
- Decompose $f=f_X+f_\perp$ with $f_\mathcal{X}\in\mathcal{H}_X$, $f_\perp\in\mathcal{H}_\perp$
- $f(x_i) = \langle f_X + f_\perp, k(x_i, \cdot)
 angle_{\mathcal{H}} = \langle f_X, k(x_i, \cdot)
 angle_{\mathcal{H}}$
- $ullet \ \|f\|^2_{\mathcal{H}} = \|f_X\|^2_{\mathcal{H}} + \|f_ot\|^2_{\mathcal{H}}$
- Minimizer needs $f_{\perp}=0$, and so $\hat{f}=\sum_{i=1}^n lpha_i k(x_i,\cdot)$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\sum_{i=1}^n \left(\sum_{j=1}^n lpha_j k(x_i,x_j)-y_i
ight)^2 = \sum_{i=1}^n \left([Klpha]_i-y_i
ight)^2$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\sum_{i=1}^n \left(\sum_{j=1}^n lpha_j k(x_i,x_j) - y_i
ight)^2 = \sum_{i=1}^n \left([Klpha]_i - y_i
ight)^2 = \|Klpha - y\|^2$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$egin{aligned} &\sum_{i=1}^n lpha_j k(x_i,x_j) - y_i \end{pmatrix}^2 &= \sum_{i=1}^n \left([Klpha]_i - y_i
ight)^2 = \|Klpha - y\|^2 \ &= lpha^\mathsf{T} K^2 lpha - 2y^\mathsf{T} K lpha + y^\mathsf{T} y \end{aligned}$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$egin{aligned} &\sum_{i=1}^n lpha_j k(x_i,x_j) - y_i \end{pmatrix}^2 &= \sum_{i=1}^n \left([Klpha]_i - y_i
ight)^2 = \|Klpha - y\|^2 \ &= lpha^\mathsf{T} K^2 lpha - 2y^\mathsf{T} K lpha + y^\mathsf{T} y \end{aligned}$$

$$\left\|\sum_{i=1}^n lpha_i k(x_i,\cdot)
ight\|^2 = \sum_{i=1}^n \sum_{j=1}^n lpha_i k(x_i,x_j) lpha_j$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$egin{aligned} &\sum_{i=1}^n lpha_j k(x_i,x_j) - y_i \end{pmatrix}^2 &= \sum_{i=1}^n \left([Klpha]_i - y_i
ight)^2 = \|Klpha - y\|^2 \ &= lpha^\mathsf{T} K^2 lpha - 2y^\mathsf{T} K lpha + y^\mathsf{T} y \end{aligned}$$

$$\left\|\sum_{i=1}^n lpha_i k(x_i,\cdot)
ight\|^2 = \sum_{i=1}^n \sum_{j=1}^n lpha_i k(x_i,x_j) lpha_j = lpha^\mathsf{T} K lpha$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$\hat{lpha} = rgmin_{lpha \in \mathbb{R}^n} lpha^\mathsf{T} K^2 lpha - 2 y^\mathsf{T} K lpha + y^\mathsf{T} y + n \lambda lpha^\mathsf{T} K lpha$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

$$egin{aligned} \hat{lpha} &= rg\min lpha^{\mathsf{T}} K^2 lpha - 2y^{\mathsf{T}} K lpha + y^{\mathsf{T}} y + n\lambda lpha^{\mathsf{T}} K lpha \ &= rg\min lpha^{\mathsf{T}} K (K + n\lambda I) lpha - 2y^{\mathsf{T}} K lpha \ &= lpha \in \mathbb{R}^n \end{aligned}$$

$$\hat{f} = rgmin_{f \in \mathcal{H}} rac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

How to find \hat{f} ? Representer Theorem: $\hat{f} = \sum_{i=1}^n \hat{lpha}_i k(x_i, \cdot)$

$$egin{aligned} \hat{lpha} &= rg\min lpha^{\mathsf{T}} K^2 lpha - 2y^{\mathsf{T}} K lpha + y^{\mathsf{T}} y + n \lambda lpha^{\mathsf{T}} K lpha \ &= rg\min lpha^{\mathsf{T}} K (K + n \lambda I) lpha - 2y^{\mathsf{T}} K lpha \ &lpha \in \mathbb{R}^n \end{aligned}$$

Setting derivative to zero gives $K(K+n\lambda I)\hat{lpha}=Ky,$ satisfied by $\hat{lpha}=(K+n\lambda I)^{-1}y$

Other kernel algorithms

• Representer theorem applies if old R strictly increasing:

$$\min_{f\in\mathcal{H}}L(f(x_1),\cdots,f(x_n))+R(\|f\|_{\mathcal{H}})$$

- Classification algorithms:
 - Support vector machines: L is hinge loss
 - Kernel logistic regression: *L* is logistic loss
- Principal component analysis, canonical correlation analysis
- Many, many more...

Some theory

- If $\mathcal H$ universal, $f\in\mathcal H$ can approximate any continuous func
 - $\int_{\mathcal{X}} \int_{\mathcal{X}} k(x,y) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) > 0$ for all nonzero finite signed measures μ
 - True for Gaussian, many other common kernels (but no finite-dimensional ones!)
 - Norm may go to ∞ as approximation gets better

Some theory

- If ${\mathcal H}$ universal, $f\in {\mathcal H}$ can approximate any continuous func
 - $\int_{\mathcal{X}} \int_{\mathcal{X}} k(x,y) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) > 0$ for all nonzero finite signed measures μ
 - True for Gaussian, many other common kernels (but no finite-dimensional ones!)
 - Norm may go to ∞ as approximation gets better
- If RKHS norm is small, can learn quickly
 - e.g. Rademacher complexity of $\{f \in \mathcal{H}: \|f\|_{\mathcal{H}} \leq B\}$

is at most
$$\sqrt{rac{B^2}{m} \sup_{x \in \mathcal{X}} k(x,x)}$$

Limitations of kernel-based learning

- Generally bad at learning *sparsity*
 - e.g. $f(x_1,\ldots,x_d)=3x_2-5x_{17}$ for large d

Limitations of kernel-based learning

- Generally bad at learning *sparsity*
 - e.g. $f(x_1,\ldots,x_d)=3x_2-5x_{17}$ for large d
- Provably slower than deep learning for a few problems
 - e.g. to learn a single ReLU, $\max(0, w^{\mathsf{T}}x)$, need norm exponential in d [Yehudai/Shamir NeurIPS-19]
 - Also some hierarchical problems, etc [Kamath+ COLT-20]

Relationship to deep learning

• Deep models usually end as $f_L(x) = w_L^\mathsf{T} f_{L-1}(x)$

Relationship to deep learning

- Deep models usually end as $f_L(x) = w_L^\mathsf{T} f_{L-1}(x)$
- Can think of as *learned* kernel, $k(x,y)=f_{L-1}(x)f_{L-1}(y)$
- Deep models usually end as $f_L(x) = w_L^\mathsf{T} f_{L-1}(x)$
- Can think of as *learned* kernel, $k(x,y) = f_{L-1}(x) f_{L-1}(y)$
- Does this gain us anything?

- Deep models usually end as $f_L(x) = w_L^\mathsf{T} f_{L-1}(x)$
- Can think of as *learned* kernel, $k(x,y) = f_{L-1}(x) f_{L-1}(y)$
- Does this gain us anything?
 - Random nets with trained last layer (NNGP) can be decent

- Deep models usually end as $f_L(x) = w_L^\mathsf{T} f_{L-1}(x)$
- Can think of as *learned* kernel, $k(x,y) = f_{L-1}(x) f_{L-1}(y)$
- Does this gain us anything?
 - Random nets with trained last layer (NNGP) can be decent
 - As width → ∞, nets become neural tangent kernel
 Widely used theoretical analysis
 - SVMs with NTK can be great on small data

- Deep models usually end as $f_L(x) = w_L^\mathsf{T} f_{L-1}(x)$
- Can think of as *learned* kernel, $k(x,y) = f_{L-1}(x) f_{L-1}(y)$
- Does this gain us anything?
 - Random nets with trained last layer (NNGP) can be decent
 - As width → ∞, nets become neural tangent kernel
 Widely used theoretical analysis
 - SVMs with NTK can be great on small data
 - Inspiration: learn the kernel model end-to-end
 - Ongoing area; good results in two-sample testing, GANs, density estimation, meta-learning, semisupervised learning, ...

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
angle_{\mathcal{H}}$

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_\mathbb P$, $\mathbb E_{X\sim\mathbb P} f(X)=\langle f,\mu_\mathbb P
 angle_{\mathcal H}$

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$, $\mathbb E_{X\sim\mathbb P} f(X)=\langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

 $\mathbb{E}_{X\sim\mathbb{P}} f(X) = \mathbb{E}_{X\sim\mathbb{P}} \langle f, k(X,\cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim\mathbb{P}} \, k(X,\cdot)
angle_{\mathcal{H}}$

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$, $\mathbb E_{X\sim\mathbb P} f(X)=\langle f,\mu_{\mathbb P}
 angle_{\mathcal H}$

 $\mathbb{E}_{X\sim\mathbb{P}} f(X) = \mathbb{E}_{X\sim\mathbb{P}} \langle f, k(X,\cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim\mathbb{P}} \ k(X,\cdot)
angle_{\mathcal{H}}$

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_\mathbb P$, $\mathbb E_{X\sim\mathbb P} f(X)=\langle f,\mu_\mathbb P
 angle_{\mathcal H}$

$$\mathbb{E}_{X\sim\mathbb{P}} f(X) = \mathbb{E}_{X\sim\mathbb{P}} \langle f, k(X,\cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim\mathbb{P}} \ k(X,\cdot)
angle_{\mathcal{H}}$$

Last step assumed e.g. $\mathbb{E}\sqrt{k(X,X)} < \infty$

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_{\mathbb P}$, $\mathbb E_{X\sim \mathbb P} f(X) = \langle f, \mu_{\mathbb P}
 angle_{\mathcal H}$

$$\mathbb{E}_{X\sim \mathbb{P}} f(X) = \mathbb{E}_{X\sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X\sim \mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

Last step assumed e.g. $\mathbb{E}\sqrt{k(X,X)} < \infty$

• Okay. Why?

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_\mathbb P$, $\mathbb E_{X\sim\mathbb P} f(X)=\langle f,\mu_\mathbb P
 angle_{\mathcal H}$

$$\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{X \sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X \sim \mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

- Last step assumed e.g. $\mathbb{E}\sqrt{k(X,X)} < \infty$
- Okay. Why?
 - One reason: ML on distributions [Szabó+ JMLR-16]

- Represent point $x \in \mathcal{X}$ as $\phi(x)$, $f(x) = \langle f, k(x, \cdot)
 angle_{\mathcal{H}}$
- Represent distribution $\mathbb P$ as $\mu_\mathbb P$, $\mathbb E_{X\sim\mathbb P}\,f(X)=\langle f,\mu_\mathbb P
 angle_{\mathcal H}$

$$\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{X \sim \mathbb{P}} \langle f, k(X, \cdot)
angle_{\mathcal{H}} = \langle f, \mathbb{E}_{X \sim \mathbb{P}} \ k(X, \cdot)
angle_{\mathcal{H}}$$

- Last step assumed e.g. $\mathbb{E} \sqrt{k(X,X)} < \infty$
- Okay. Why?
 - One reason: ML on distributions [Szabó+ JMLR-16]
 - More common reason: comparing distributions

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

• Last line is Integral Probability Metric (IPM) form

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}}\,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}}\,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \: k(t, X) - \mathbb{E}_{\mathbb{Q}} \: k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \, k(t, X) - \mathbb{E}_{\mathbb{Q}} \, k(t, Y)$$



$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \, k(t, X) - \mathbb{E}_{\mathbb{Q}} \, k(t, Y)$$



$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot) \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} k(t, X) - \mathbb{E}_{\mathbb{Q}} k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \, k(t, X) - \mathbb{E}_{\mathbb{Q}} \, k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \, k(t, X) - \mathbb{E}_{\mathbb{Q}} \, k(t, Y)$$

$$egin{aligned} \mathrm{MMD}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}\|_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \langle f,\mu_{\mathbb{P}}-\mu_{\mathbb{Q}}
angle_{\mathcal{H}} \ &= \sup_{\|f\|_{\mathcal{H}}\leq 1} \mathbb{E}_{X\sim\mathbb{P}} \,f(X) - \mathbb{E}_{Y\sim\mathbb{Q}} \,f(Y) \end{aligned}$$

- Last line is Integral Probability Metric (IPM) form
- f is called "witness function" or "critic": high on $\mathbb P$, low on $\mathbb Q$

$$f^*(t) \propto \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, k(t, \cdot)
angle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}} \, k(t, X) - \mathbb{E}_{\mathbb{Q}} \, k(t, Y)$$

More!

- Foundations: Berlinet and Thomas-Agnan, *RKHS in Probability and Statistics*
- Hardcore theoretical details: Steinwart and Christmann, *Support Vector Machines*
- Close connections to Gaussian processes [Kanagawa+ 'GPs and Kernel Methods' 2018]
- Mean embeddings: survey of [Muandet+ 'Kernel Mean Embedding of Distributions']
- The practical sessions! Some pointers in there