# Kernel Methods Notes

- A **Kernel** is a function $K: \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ such that there exists a Hilbert space $\mathcal{H}$ and mapping $\phi: \mathcal{X} \rightarrow \mathcal{H}$ where $K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$

- A **Hilbert Space** is a vector space on which an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \longrightarrow \mathbb{R}$ is defined, this having the following properties:

  - $\langle a f_1 + b f_2, g \rangle_{\mathcal{H}} = a \langle f_1, g \rangle_{\mathcal{H}} + b \langle f_2, g \rangle_{\mathcal{H}}$ (linear)

  - $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$ (symmetric)

  - $\langle f, f \rangle_{\mathcal{H}} \geq 0$, $= 0$ only when $f = 0$

- All kernels $K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$ are __positive definite__ functions:

  given arbitrary $a_1, \ldots, a_n \in \mathbb{R}$, $x_1, \ldots, x_n \in \mathcal{X}$

$$\sum_i \sum_j a_i a_j K(x_i, x_j) = \sum_i \sum_j \langle a_i \phi(x_i), a_j \phi(x_j) \rangle_{\mathcal{H}}$$

$$= \left\langle \sum_i a_i \phi(x_i), \sum_j a_j \phi(x_j) \right\rangle_{\mathcal{H}}$$

$$= \left\| \sum_i a_i \phi(x_i) \right\|_{\mathcal{H}}^2 \geq 0 \quad \dashv$$

- It turns out that the opposite direction holds as well: all positive definite functions are kernels!
- Therefore, all sums of kernels $K(x, x') = K_1(x, x') + K_2$ are kernels! for arbitrary $a_1, \dots, a_n \in \mathbb{R}$, $x_1, \dots, x_n \in \mathcal{X}$

$$\sum_i \sum_j a_i a_j K(x_i, x_j) = \sum_i \sum_j a_i a_j \left( K_1(x_i, x_j) + K_2(x_i, x_j) \right)$$

$$= \left\| \sum_i a_i \phi_1(x_i) \right\|_{\mathcal{H}_1}^2 + \left\| \sum_i a_i \phi_2(x_i) \right\|_{\mathcal{H}_2}^2 \geq$$

$$\Rightarrow \text{positive-definite} \therefore \text{a kernel}$$

- All products of kernels $K(x, x') = K_1(x, x') K_2(x, x')$ are kernels!

$$K_1(x, x') K_2(x, x') = \langle \phi_1(x), \phi_1(x') \rangle_{\mathcal{H}_1} \langle \phi_2(x), \phi_2(x') \rangle_{\mathcal{H}_2}$$

can always take trace of a scalar $\left( \quad = \phi_1(x')^T \phi_1(x) \, \phi_2(x)^T \phi_2(x') \right.$

$$= \phi_1(x')^T \phi_1(x^*) \, \text{Trace} \left[ \phi_2(x') \phi_2(x)^T \right]$$

can move a scalar into a trace $\left( \quad = \text{Tr} \left[ \phi_2(x') \underbrace{\phi_1(x')^T \phi_1(x^*)}_{A^T} \underbrace{\phi_2(x)^T}_{B} \right] \right.$

Frobenius product $\left( \quad \overset{*}{=} \text{Tr} \left[ A^T B \right] \right.$

$$= \text{vec}(A)^T \text{vec}(B)$$

$$= \left\langle \text{vec}\left( \phi_2(x') \phi_2(x')^T \right), \text{vec}\left( \phi_1(x^*) \phi_2(x)^T \right) \right\rangle_{\mathcal{H}}$$

$$= \langle \psi(x'), \psi(x) \rangle_{\mathcal{H}} = K(x, x') \checkmark$$

- Every kernel is associated with a unique RKHS $\mathcal{H}$, which has the following properties:

  - $\forall x \in \mathcal{X}, \; K(\cdot, x) \in \mathcal{H}$.
  - $\forall x \in \mathcal{X}, \; \forall f \in \mathcal{H}, \quad \langle f, K(\cdot, x) \rangle = f(x)$

  $$\underline{\text{reproducing property}}$$

- Ex. RKHS defined by a Fourier series

consider the space of all periodic functions on $[-\pi, \pi]$:
$$f(x) = \sum_{\ell = -\infty}^{\infty} \hat{f}_\ell \, e^{i\ell x}$$

We can then define the $\infty - D$ space spanned by the orthonormal basis $\{e^{i\ell x}\}_{\ell = -\infty}^{\infty}, \; x \in \mathbb{R}$ together with the standard $L2$ dot product $\langle \cdot, \cdot \rangle$, to give us a Hilbert space $\mathcal{H}$, where $\langle f, g \rangle_{L2} = \sum_{\ell = -\infty}^{\infty} \hat{f}_\ell \overline{\hat{g}_\ell}$.
Is $\mathcal{H}$ an RKHS? Let $K(x, y) = K(x - y)$
we check for the reproducing property:
$$= \sum_{\ell = -\infty}^{\infty} \hat{K}_\ell \, e^{i\ell(x-y)}$$

$$\langle f, K(\cdot, x) \rangle_{L2} = \sum_{\ell = -\infty}^{\infty} \hat{f}_\ell \, \overline{\hat{K}_\ell \, e^{i\ell y}}$$

$$= \sum_{\ell = -\infty}^{\infty} \hat{K}_\ell \, e^{i\ell y} \, e^{i\ell}$$

$$= \sum_{\ell = -\infty}^{\infty} \hat{K}_\ell \, \hat{f}_\ell \, e^{i\ell x} \; \neq \; f(x)$$

(Given this kernel, what is the dot product of the associated RKHS?)

So $\mathcal{H}$ is $\underline{\text{not}}$ an RKHS. But we can easily modify it so that it is: $\mathcal{H}^*$ with $\langle f, g \rangle_{\mathcal{H}^*} = \sum_{\ell = -\infty}^{\infty} \dfrac{\hat{f}_\ell \, \overline{\hat{g}_\ell}}{\hat{K}_\ell}$

Now, $\left\langle f, K(\cdot, x)\right\rangle_{\mathcal{H}^*} = \sum\limits_{\ell=-\infty}^{\infty} \dfrac{\hat{f_\ell} \hat{K_\ell} e^{i\ell x}}{\hat{K_\ell}} = \sum\limits_{\ell=-\infty}^{\infty} \hat{f_\ell} e^{i\ell x} = f(x)$

$\left\langle K(\cdot, x), K(\cdot, y)\right\rangle_{\mathcal{H}^*} = \sum\limits_{\ell=-\infty}^{\infty} \dfrac{\hat{K_\ell} e^{-i\ell x} \hat{K_\ell} e^{i\ell y}}{\hat{K_\ell}} = \sum\limits_{\ell=-\infty}^{\infty} \hat{K_\ell} e^{i\ell(y-x)} = K(y-x)$

Importantly, $\left\langle f, f\right\rangle_{\mathcal{H}^*} = \|f\|_{\mathcal{H}^*}^2 = \sum\limits_{\ell=-\infty}^{\infty} \dfrac{|\hat{f_\ell}|^2}{\hat{K_\ell}}$, so the kernel enforces smoothness since any $f \in \mathcal{H}^*$ must have $\hat{f_\ell}$ that decay faster than $\hat{K_\ell}$ for $\|f\|_{\mathcal{H}^*}^2 < \infty$, i.e. $f(\cdot)$ must be at least as smooth (low amplitudes at higher frequencies) as $K(\cdot)$.

— **Kernel PCA**: just like normal PCA but performed in feature space, via the reproducing property:

$f^* = \underset{\substack{f \in \mathcal{H} \\ \|f\|_{\mathcal{H}}=1}}{\arg\max}$ variance of data projected into $\mathcal{H}$ via feature map $\phi(x) = K(x, \cdot)$ along unit vector $f$

$= \underset{\|f\|_{\mathcal{H}}=1}{\arg\max} \left\langle f, \text{~~~~~~~~~~~~} \right\rangle_{\mathcal{H}} \qquad C = \dfrac{1}{N}\sum\limits_{i}\left(\phi(x_i) - \dfrac{1}{N}\sum\limits_{j}\phi(x_j)\right)^2$

$= \underset{\|f\|_{\mathcal{H}}=1}{\arg\max} \dfrac{1}{N}\sum\limits_{i=1}^{N}\left(\left\langle f, \phi(x_i)\right\rangle - \bar{\phi}\right)_{\mathcal{H}}^2 = \dfrac{1}{N}\sum\limits_{i}^{N}\tilde{\phi}(x_i)^2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \bar{\phi} = \dfrac{1}{N}\sum\phi(x_i)$

$\dfrac{1}{N}\sum\limits_{i}\left\langle f, \tilde{\phi}(x_i)\right\rangle\left\langle f, \tilde{\phi}(x_i)\right\rangle \qquad \tilde{\phi}(x_i) = \phi(x_i) - \bar{\phi}$

$\dfrac{1}{N}\sum\limits_{i}\left\langle f, \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i) f\right\rangle$

$= \underset{\|f\|_{\mathcal{H}}}{\arg\max} \left\langle f, C f\right\rangle, \qquad C = \dfrac{1}{N}\sum\limits_{i=1}^{N} \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i)$

$$\Rightarrow \frac{\partial}{\partial \delta}\left[\langle \delta, C\delta \rangle_{\mathcal{H}} + \lambda\left(\langle \delta,\delta\rangle_{\mathcal{H}} - 1\right)\right] = 0$$

$$\iff C\delta = \lambda\delta$$

$$\Rightarrow \delta^* = \text{largest e-vector } C$$

but this requires computing $C$, which $\not{p}$
lives $\text{\sout{neural}}$ in $\mathbb{R}^{\infty \times \infty}$
$\rightarrow$ How can we avoid feature space?

$\Rightarrow$ We can always express $f$ as a linear combination of data points, without loss of generality, since any dimensions orthogonal to the space spanned by $\{\tilde{\phi}(x_i)\}_{i=1}^{n}$ will disappear in the first like $\langle \delta, \tilde{\phi}(x_i)\rangle_{\mathcal{H}}$, thus rendering them irrelevant to the optimization:

$$f = \sum_{i=1}^{N} \alpha_i \tilde{\phi}(x_i)$$

$$\tilde{K}(x,x') = \langle \tilde{\phi}(x), \tilde{\phi}(x')\rangle_{\mathcal{H}}$$

$$\iff f(\cdot) = \sum_{i=1}^{N} \alpha_i \tilde{K}(x_i, \cdot) \qquad \left(\text{by reproducing property}\right)$$

Thus we need only solve for the $\alpha$'s:

$$Cf = \frac{1}{N}\sum_{i=1}^{N} \tilde{\phi}(x_i)\sum_{j=1}^{N}\alpha_j\langle \tilde{\phi}(x_i), \tilde{\phi}(x_j)\rangle_{\mathcal{H}}$$

$$\alpha_j \tilde{K}(x_i)$$

$$= \frac{1}{N}\sum_{i=1}^{N} \tilde{\phi}(x_i)\sum_{j=1}^{N}\alpha_j\tilde{K}(x_i,x_j) \Rightarrow \langle \tilde{\phi}(x_k), Cf\rangle = \frac{1}{N}\sum_{i}\tilde{K}(x_k,x_i)$$

$$\langle \tilde{\phi}(x_k), \lambda f\rangle_{\mathcal{H}} = \lambda\sum_{i}\alpha_i\tilde{K}(x_k,x_i) \Rightarrow \frac{1}{N}\tilde{K}\tilde{K}\alpha = \lambda\tilde{K}\alpha$$

where $\tilde{K}_{ij} = \tilde{K}(x_i, x_j)$. Since this matrix is symmetric and positive semidefinite, its inverse exists, so we get the following eigenvalue equation:

$$\tilde{K}\alpha = N\lambda\alpha$$

So we can solve for $\alpha$ by constructing the Gram matrix $\tilde{K}$ and solving the eigenvalue equation, giving us the directions of $\phi$ of greatest variance without having to work out all in feature space. (ie. biggest ...

Importantly, $\phi$ is a function, so kernel PCA, as opposed to regular PCA, can give us ......... non-linear principal subspaces rather than just ............... hyperplanes (depending on the kernel).

— <u>Kernel Ridge Regression</u>: ridge regression in feature space

$$y = \phantom{xxxx} w^T \phi(x) + \epsilon, \quad \phi(x) \in \mathcal{H}$$

$$\Rightarrow w^* = \underset{w \in \mathcal{H}}{\arg\min} \left[ \sum_{i=1}^{N} \left( y_i - \langle w, \phi(x_i) \rangle_{\mathcal{H}} \right)^2 + \lambda \|w\|_{\mathcal{H}}^2 \right]$$

$$= \underset{w \in \mathcal{H}}{\arg\min} \left[ \|Y - X^T w\|_{\mathcal{H}}^2 + \lambda \|w\|_{\mathcal{H}}^2 \right], \quad X = \left[ \phi(x_1) \cdots \phi(x_i) \right]$$

$$= \underset{w \in \mathcal{H}}{\arg\min} \left[ Y^T Y - 2 Y^T X^T w + w^T \left( X X^T + \lambda I \right) w \right]$$

completing the square

$$= \underset{w \in \mathcal{H}}{\arg\min} \left[ Y^T Y + \left\| (X X^T + \lambda I)^{\frac{1}{2}} w - (X X^T + \lambda I)^{-\frac{1}{2}} X Y \right\|_{\mathcal{H}}^2 - \left\| (X X^T + \lambda I)^{-\frac{1}{2}} X \cdots \right. \right.$$

$$= \left( X X^T + \lambda I \right)^{-1} X Y$$

(we could've done this by taking derivatives, but derivatives don't necessarily exist for discrete $x_i, y_i$,)

To avoid having to do anything in feature space, we rewrite this in terms of the (gram) matrix $K = X^T X$:

① Via SVD:

$K_{ij} = K(x_i, x_j)$

$$\underset{D \times N}{X} = \underset{D \times D}{\begin{bmatrix} \tilde{U} \end{bmatrix}} \underset{D \times N}{\begin{bmatrix} \tilde{S} \\ 0 \end{bmatrix}} \underset{N \times N}{\begin{bmatrix} \tilde{V} \end{bmatrix}}$$

$\underbrace{\phantom{XXXX}}_{(\text{orthogonal})} \underbrace{\phantom{XXXX}}_{(\text{diagonal})} \underbrace{\phantom{XXXX}}_{(\text{orthogonal})}$

Let
$$U = \tilde{U} \qquad D \times D$$
$$S = \begin{bmatrix} \tilde{S} & 0 \\ 0 & 0 \end{bmatrix} \qquad D \times D$$
$$V = \begin{bmatrix} \tilde{V} & 0 \end{bmatrix} \qquad N \times D$$

such that $X = U S V^T$

we then have:

$$w^* = \left( U S^2 U^T + \lambda I \right)^{-1} U S V^T Y$$

$$= U \left( S^2 + \lambda I \right)^{-1} U^T U S V^T Y$$

can do this since $S$ is diagonal and square (hence the change from the $\tilde{\phantom{x}}$ using SVD)

$$= U S \left( S^2 + \lambda I \right)^{-1} V^T Y$$

$$= U S V^T V \left( S^2 + \lambda I \right)^{-1} V^T Y$$

$$= U S V^T \left( V^T S^2 V + \lambda I \right)^{-1} Y$$

$$= X \left( X^T X + \lambda I \right)^{-1} Y$$

$$= \underline{\underline{X \left( K + \lambda I \right)^{-1} Y}}$$

② Via Woodberry Identity:

$$w^* = (XX^T + \lambda I)^{-1} XY$$

$$= \left(\lambda^{-1} I - \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1} X^T \lambda^{-1}\right) XY$$

$$= \left[\lambda^{-1} X - \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1} \lambda^{-1} X^T X\right] Y$$

$$= \left[\lambda^{-1} X + \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1}\right.$$
$$- \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1}$$
$$\left. - \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1} \lambda X^T X\right] Y$$

$$= \left[\lambda^{-1} X + \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1}\right.$$
$$\left. - \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1} (\lambda^{-1} X^T X + I)\right]$$

$$= \lambda^{-1} X (\lambda^{-1} X^T X + I)^{-1} Y$$

$$= \underbrace{X (X^T X + \lambda I)^{-1}}_{K} Y$$

Thus, our optimal weights are a weighted sum of the data points: $w^* = \sum_i \alpha_i \phi(x_i)$, $\underline{\alpha} = (K + \lambda I)^{-1} Y$

Note that $w^*$ is a function in $\mathcal{H}$, such that its smoothness is constrained by the kernel since $\|w^*\|_{\mathcal{H}}^2 < \infty$. The larger our regularizing the constant $\lambda$, the smoother our resulting regression function $\langle w^*, \phi(x) \rangle_{\mathcal{H}} = w^*(x)$ will be

# Part II : MMD, HSIC, COCO

- Just like the "kernel trick" allows us to express functions in terms of feature space:
$$f(x) = \langle f, K(x, \cdot) \rangle_{\mathcal{H}}$$

the "mean trick" allows us to do the same with expectations:
$$\mathbb{E}_{x \sim p}[f(x)] = \langle \mu_p, f \rangle_{\mathcal{H}}$$

By the reproducing property,

$\boxed{\begin{array}{c} \text{probability} \\ \text{feature map} \end{array}}$ (analog to feature map $f(x)$) $\mu_p(x) = \langle \mu_p, K(x, \cdot) \rangle_{\mathcal{H}} = \mathbb{E}_{x \sim p}[K(x, X)]$

so we can estimate it empirically just like ~~oooo~~ usual!

$$\hat{\mu}_p(a) = \frac{1}{N} \sum_i \langle K(x_i, \cdot), K(a, \cdot) \rangle_{\mathcal{H}}$$

$$= \frac{1}{N} \sum_i K(x_i, a)$$

\qquad mean embedding

· We can prove that $\mu_p$ exists in feature space (ie. prove that the "mean trick" works) via the Riesz representation theorem:

any **bounded linear operator** $A : \mathcal{H} \to \mathbb{R}$ i.e.
$$|Af| \leq \lambda_A \|f\|_{\mathcal{H}} \quad \forall f \in \mathcal{H}$$
can be expressed as
$$Af = \langle f, g_A \rangle_{\mathcal{H}} \quad \text{for some } g \in \mathcal{H}$$

Thus, if we prove that the expectation operation $\mathbb{E}_p$ is bounded, then $\mu_p \in \mathcal{H}$:

assuming $\mathbb{E}_p\left[\sqrt{K(x,x)}\right] < \infty$

$$\left|\mathbb{E}_p f(x)\right| \overset{\text{Jensen}}{\leq} \mathbb{E}_p |f(x)| = \mathbb{E}_p\left[\langle f, K(x,\cdot)\rangle_{\mathcal{H}}\right] \overset{\substack{\text{Cauchy-}\\\text{Schwartz}}}{\leq} \mathbb{E}_p\left[\|K(x,\cdot)\|_{\mathcal{H}} \|f\|_{\mathcal{H}}\right]$$

$$\therefore \mathbb{E}_p f(x) = \langle f, \mu_p\rangle_{\mathcal{H}}, \quad \mu_p \in \mathcal{H}$$

$$= \mathbb{E}_p\left[\sqrt{\langle K(x,\cdot), K(x,\cdot)\rangle}\right]\|f\|_{\mathcal{H}}$$

$$= \mathbb{E}_p\left[\sqrt{K(x,x)}\right]\|f\|_{\mathcal{H}}$$

$$= \lambda\|f\|_{\mathcal{H}}$$

- To compare means, we use the Max. Mean Discrepancy

$$MMD(P,Q; \mathcal{H}) = \sup_{f\in\mathcal{H}}\left[\mathbb{E}_p f(x) - \mathbb{E}_Q f(y)\right] \quad \substack{\text{for } x\sim P\\ y\sim Q}$$

where $\mathcal{H}$ is the unit ball in RKHS, i.e. $\|f\|_{\mathcal{H}}\leq 1$
$$= \sup_{f\in\mathcal{H}}\langle f, \mu_p - \mu_Q\rangle_{\mathcal{H}}$$

$$= \|\mu_p - \mu_Q\|_{\mathcal{H}}$$

which we estimate empirically by:

$$MMD^2 = \langle \mu_p, \mu_p\rangle_{\mathcal{H}} + \langle \mu_Q, \mu_Q\rangle_{\mathcal{H}} - 2\langle \mu_p, \mu_Q\rangle_{\mathcal{H}}$$

$$\simeq \frac{1}{N(N-1)}\sum_{i=1}^{N}\sum_{j\neq i} K(x_i, x_j) + \frac{1}{M(M-1)}\sum_{i=1}^{M}\sum_{j\neq i} K(y_i, y_j)$$

$$- \frac{2}{NM}\sum_{i=1}^{N}\sum_{j=1}^{M} K(x_i, y_j)$$

exclude repeating the same
data point to ensure
estimator remains unbiased

- MMD $= 0$ iff $P = Q$ whenever $\mathcal{H}$ is a characteristic RKHS with characteristic kernel $K(\cdot, \cdot)$.

**Prop** For periodic on $[-\pi, \pi]$ and translation-invariant $K(\cdot, \cdot)$, $K$ is characteristic iff $\hat{K}_\ell \neq 0 \ \forall \ell$

$$\mu_P(z) = \langle \mu_P, K(z, \cdot) \rangle$$
$$= \mathbb{E}_{x \sim P} \, K(z, x)$$
$$= \mathbb{E}_{x \sim P} \, K(z - x)$$
$$= \int_{-\pi}^{\pi} K(z - x) \, dP(x)$$

$$\Rightarrow \hat{\mu}_{P,\ell} = \int \mu_P(z) \, e^{-i\ell z} \, dz$$

$$= \int_{-\pi}^{\pi} \int K(z - x) \, e^{-i\ell z} \, dP(x) \, dz$$

$$= \int_{-\pi}^{\pi} \int K(v) \, e^{-i\ell(v + x)} \, dP(x) \, dv$$

① ~~...~~

② MMD penalizes some freqs more than others, depending on Kernel smooth

$$= \int_{-\pi}^{\pi} K(v) \, e^{-i\ell v} \int e^{-i\ell x} \, dP(x)$$

*Recalling that for periodic & translation invariant $K$, $\|S\|_{\mathcal{H}}^2 = \sum_\ell \frac{|\hat{S}_\ell|^2}{\hat{K}_\ell}$*

Fourier coefs of the probability distribution p density function

$$= \hat{K}_\ell \cdot \hat{\varphi}_{P,\ell}$$

$$\Rightarrow MMD^2 = \|\mu_P - \mu_Q\|_{\mathcal{H}}^2 = \sum_{\ell = -\infty}^{\infty} \frac{|\hat{K}_\ell \hat{\varphi}_{P,\ell} - \hat{K}_\ell \hat{\varphi}_{Q,\ell}|^2}{\hat{K}_\ell} = \sum_{\ell = -\infty}^{\infty} \hat{K}_\ell \, |\hat{\varphi}_{P,\ell} - \hat{\varphi}_{Q,\ell}|^2$$

For any function on $\mathbb{R}^D$ we can show the following, via Bochner's theorem:

$$MMD^2 = \frac{1}{2} \int_{\mathbb{R}^D} |\hat{\varphi}_P(l) - \hat{\varphi}_Q(l)|^2 \, d\Lambda(l)$$

Fourier transform of $K$

which says the same thing.

$\Rightarrow$ $K$ characteristic iff $\text{supp}(\Lambda) = \mathbb{R}^D$

$\Rightarrow$ any continuous $K$ with ~~addd~~ Fourier transform $\Lambda$ s.t. $\text{supp}(\Lambda) = \mathbb{R}^D$ is characteristic

$\hookrightarrow$ support = set of dev that are **not** mapped to 0

— For hypothesis testing, use
$$\widehat{MMD}^2 = \frac{1}{N(N-1)} \sum^N \sum_{i \neq i} K(x_i, x_j) + K(y_i, y_j) - K(x_i, \cdot) - K(y_i, \cdot)$$
(lower variance since dropped some terms, but still unbiased) ~~$x_i y_j$~~

$H_1 : P \neq Q$

$$\sqrt{N}\left(\widehat{MMD}^2 - MMD^2\right) \sim \mathcal{N}(0, \sigma_u^2)$$

$\hookrightarrow$ asymptotically normal

$H_0 : P = Q$

$$N \cdot \widehat{MMD} \sim \sum_{l \geq 1}^{\infty} \lambda_l (z_l^2 - 2), \quad z_l \overset{i.i.d.}{\sim} \mathcal{N}(0, 2)$$

$\hookrightarrow$ degenerate $U$-statistic, so need to estimate via e.g. permutation
Pearson moment matching

- Just like we can show that $\mathbb{E}_p[\delta(x)]$ can be expressed in feature space via the mean embedding $\mu$, we can show that the cross-covariance

$$\mathbb{E}_{P_{xy}}[f(x,y)] = \mathbb{E}_{P_{xy}}\left[\left\langle \phi(x) \otimes \psi(y), f \right\rangle_{\mathcal{F} \times G}\right]$$

$$= \left\langle \tilde{C}_{xy}, f \right\rangle_{\mathcal{F} \times G}$$

for feature maps $\phi : \mathcal{X} \to \mathcal{F}$, $\psi : \mathcal{Y} \to G$, and Hilbert–Schmidt operators $f, \tilde{C}_{xy} \in \mathcal{F} \times G$

We can again show $\tilde{C}_{xy}$ exists via Riesz representer theorem:

$$\left| \mathbb{E}_{P_{xy}}[f(x,y)] \right| \leq \mathbb{E}_{P_{xy}}\left| f(x,y) \right| = \mathbb{E}_{P_{xy}}\left| \left\langle f, \phi(x) \otimes \psi(y) \right\rangle \right|$$

Jensen

$$\overset{\text{Cauchy-}}{\underset{\text{Schwartz}}{\leq}} \mathbb{E}_{P_{xy}}\left[ \|f\|_{\mathcal{F} \times G} \|\phi(x) \otimes \psi(y)\|_{\mathcal{F}} \right]$$

see tensor product norm

$$= \|f\|_{\mathcal{F} \times G} \mathbb{E}_{P_{xy}}\left[\|\phi(x)\|_{\mathcal{F}} \|\psi\|\right]$$

Hence by Riesz, the bounded linear operator
$\mathbb{E}_{P_{xy}}[f(x,y)]$ can be expressed

$$= \|f\|_{\mathcal{F} \times G} \mathbb{E}_{P_{xy}}\left[\sqrt{k(x,x)} \, l(y,y)\right]$$

as $\left\langle \tilde{C}_{xy}, f \right\rangle_{\mathcal{F} \times G}$

$$< \infty$$

We can see that $\tilde{C}_{xy}$ gives us the cross-covariance b/w variables in feature space by considering

$$\mathbb{E}_{P_{xy}}\left[k(x,X) \, l(y,Y)\right] = \mathbb{E}_{P_{xy}}\left[\left\langle k(x,\cdot), k(X,\cdot)\right\rangle_{\mathcal{F}} \left\langle l(y,\cdot), l(Y,\cdot)\right\rangle_{G}\right]$$

$$= \mathbb{E}_{P_{xy}}\left[\left\langle k(x,\cdot), k(X,\cdot) \otimes l(Y,\cdot) \, l(y,\cdot)\right\rangle_{\mathcal{F}}\right]$$

$$= \left\langle k(x,\cdot), \underset{P_{xy}}{\mathbb{E}} \left[\phi(x) \otimes \psi(y)\right] l(y,\cdot) \right\rangle_{\mathcal{F}}$$

$$= \left\langle k(x,\cdot) \otimes l(y,\cdot), \underset{P_{xy}}{\mathbb{E}} \phi(x) \otimes \psi(y) \right\rangle_{\mathcal{F} \times \mathcal{G}}$$

$$= \left\langle k(x,\cdot) \otimes l(y,\cdot), \tilde{C}_{xy} \right\rangle_{\mathcal{F} \times \mathcal{G}}$$

where $k(x,\cdot)$, $l(y,\cdot)$ are the feature maps of
two variables $x \in \mathcal{X}$, $y \in \mathcal{Y}$.
The centered cross-covariance is then

$$C_{xy} = \underset{P_{xy}}{\mathbb{E}}\left[\phi(x) \otimes \psi(y)\right] - \underset{P_x}{\mathbb{E}}\left[\phi(x)\right] \underset{P_y}{\mathbb{E}}\left[\psi\right.$$

$$= \tilde{C}_{xy} - \mu_x \overset{\otimes}{} \mu_y \qquad \otimes$$

Which we can estimate empirically by:

$$\hat{C}_{xy} := \frac{1}{N} \sum_{i=1}^{N} \phi(x_i) \otimes \psi(y_i) - \hat{\mu}_x \otimes \hat{\mu}_y, \quad \hat{\mu}_x = \frac{1}{N}\sum_i \phi(x_i)$$

We can write this in matrix notation using the
centering matrix $H = I_{N \times N} - \frac{1}{N} \mathbb{1}_{N \times N}$:

$$\boxed{\hat{C}_{xy} = \frac{1}{N} X H Y^T} = \frac{1}{N} \tilde{X} \tilde{Y}^T = \frac{1}{N} \sum_i (\phi(x_i) - \hat{\mu}_x) \otimes \psi(y_i)$$

$$= \frac{1}{N} \sum_i \phi(x_i) \otimes \psi(y_i) - \hat{\mu}_x \otimes \frac{1}{N}\sum_i$$

$$= \frac{1}{N} \sum_i \phi(x_i) \otimes \psi(y_i) - \hat{\mu}_x \hat{\mu}_y$$

$$X = \left[\phi(x_1) \cdots \phi(x_N)\right], \quad Y = \left[\psi(y_1) \cdots \psi(y_N)\right]$$

- <u>Hilbert-Schmidt operators</u> are like matrices in $\mathcal{F} \times \mathcal{G}$, e.g. Let $L, M \in \mathcal{F} \times \mathcal{G}$ s.t. $L, M : \mathcal{G} \to \mathcal{F}$. Suppose $\{f_i\}_{i \in I}, \{g_j\}_{j \in J}$ are bases for $\mathcal{F}$ and $\mathcal{G}$, respectively.

We then define the <u>HS norm</u>:

$$\|L\|_{HS}^2 := \sum_{j \in J} \|L g_j\|_{\mathcal{F}}^2$$

$$= \sum_{i \in I} \sum_{j \in J} \left| \langle L g_j, f_i \rangle \right|^2$$

The HS inner product is then:

$$\langle L, M \rangle_{HS} = \sum_{j \in J} \langle L g_j, M g_j \rangle_{\mathcal{F}}$$

$$= \sum_{i \in I} \sum_{j \in J} \langle L g_j, f_i \rangle_{\mathcal{F}} \langle M g_j, f_i \rangle_{\mathcal{F}}$$

For a rank 1 operator $a \otimes b$, we have:

$$\|a \otimes b\|_{HS}^2 = \sum_{j \in J} \|a \otimes b \, g_j\|_{\mathcal{F}}^2 \qquad a \in \mathcal{F}, b \in \mathcal{G}$$

$$\cancel{\sum_{j \in J} \cdots}$$

$$= \sum_{j \in J} \|a \langle b, g_j \rangle_{\mathcal{G}}\|_{\mathcal{F}}^2 = \|a\|_{\mathcal{F}}^2 \sum_{j \in J} \langle b, g_j \rangle_{\mathcal{G}}^2$$

$$= \|a\|_{\mathcal{F}}^2 \|b\|_{\mathcal{G}}^2$$

$$\langle L, a \otimes b \rangle_{HS} = \sum_{j \in S} \langle L g_j, a \otimes b g_j \rangle_E$$

$$= \sum_{j \in S} \langle L g_j, a \rangle_E \langle b, g_j \rangle_G$$

$$= \left\langle \sum_{j \in S} L \langle b, g_j \rangle g_j, a \right\rangle_E$$

$$= \langle L b, a \rangle_E$$

$$\langle u \otimes v, a \otimes b \rangle_{HS} = \langle (u \otimes v) b, a \rangle_E$$

$$= \langle v, b \rangle_G \langle u, a \rangle_E$$

— Given $X \sim P_x$, $Y \sim P_y$ we can test for independence via <u>Hilbert-Schmidt Independence Criterion</u>

$$\text{HSIC}(P_{xy}, P_x P_y) = \text{MMD}^2(P_{xy}, P_x P_y ; \mathcal{F} \times \mathcal{G})$$

$$= \sup_{\|f\| \leq 1} \left\langle f, \mu_{P_{xy}} - \mu_{P_x} \mu_{P_y} \right\rangle_{\mathcal{F} \times \mathcal{G}}^2$$

$$= \left\| \mu_{P_{xy}} - \mu_{P_x} \mu_{P_y} \right\|_{\mathcal{F} \times \mathcal{G}}^2$$

where $f(x,y)$ is a "matrix" in $\mathcal{F} \times \mathcal{G}$ Hilbert space
and

$$\mu_{P_{xy}} = \mathbb{E}_{P_{xy}} \phi(x) \otimes \psi(y) = \tilde{C}_{xy} \in \mathcal{F} \times \mathcal{G}$$

$$\Rightarrow \mu_{P_{xy}}(x,y) = \mathbb{E}_{P_{xy}} \left[ k(x, X) \, l(y, Y) \right] \quad \left( \text{see above part on } \tilde{C}_{xy} \right.$$

$$\mu_{P_x P_y} = \mathbb{E}_{P_x} \mathbb{E}_{P_y} \phi(x) \otimes \psi(y)$$

$$= \mathbb{E}_{P_x} \phi(x) \otimes \mathbb{E}_{P_y} \psi(y)$$

$$= \mu_x \otimes \mu_y$$

Thus, we can estimate it empirically via:

$$\text{HSIC} \simeq \left\| \hat{C}_{xy} \right\|^2 = \left\| X H Y^T \right\|^2$$

$$= \text{Tr}\left[ Y H X^T X H Y^T \right]$$

$$= \text{Tr}\left[ X^T X H Y^T Y H \right]$$

$$= \text{Tr}\left[ K H L H \right]$$

But note this estimate is biased, since we are ~~including~~ including in our product the kernels evaluated at single points ⚡

$$\left(K_{ii} = K(x_i, x_i), \quad L_{ii} = \ell(y_i, y_i)\right):$$

$$HSIC = \|C_{xy}\|^2 = \|\tilde{C}_{xy} - \mu_x \otimes \mu_y\|^2$$

$$= \left\langle \tilde{C}_{xy}, \tilde{C}_{xy} \right\rangle_{\mathcal{F} \times \mathcal{G}} + \left\langle \mu_x \otimes \mu_y, \mu_x \otimes \mu_y \right\rangle_{\mathcal{F} \times \mathcal{G}}$$

$$- 2 \left\langle \tilde{C}_{xy}, \mu_x \otimes \mu_y \right\rangle_{\mathcal{F} \times \mathcal{G}}$$

$$= \left\langle \tilde{C}_{xy}, \tilde{C}_{xy} \right\rangle_{\mathcal{F} \times \mathcal{G}} + \left\langle \mu_x, \mu_x \right\rangle_{\mathcal{F}} \left\langle \mu_y, \mu_y \right\rangle_{\mathcal{G}}$$

$$- 2 E_{x,y \sim P_{xy}} \left\langle \phi(x) \otimes \psi(y), \mu_x \otimes \mu_y \right\rangle_{\mathcal{F} \times \mathcal{G}}$$

(margin note: these should really all be functions $K(x_i,), \ell(y_i,)$)

$$= E_{x,y \sim P_{xy}} E_{x',y' \sim P_{xy}} \left\langle \phi(x) \otimes \psi(y), \phi(x') \otimes \psi(y') \right\rangle_{\mathcal{F} \times \mathcal{G}}$$

$$+ E_{x \sim P_x} E_{x' \sim P_x} \left\langle \phi(x), \phi(x') \right\rangle_{\mathcal{F}} E_{y \sim P_y} E_{y' \sim P_y} \left\langle \psi(y), \psi(y') \right\rangle$$

$$- 2 E_{x,y \sim P_{xy}} E_{x' \sim P_x} E_{y' \sim P_y} \left\langle \phi(x), \phi(x') \right\rangle_{\mathcal{F}} \left\langle \psi(y), \psi(y') \right\rangle$$

$$= E_{x,y \sim P_{xy}} E_{x',y' \sim P_{xy}} K(x,x') \ell(y,y') + E_{x \sim P_x} E_{x' \sim P_x} K(x,x') E_{y \sim P_y} E_{y' \sim P_y} \ell(y,y')$$

$$- 2 E_{x,y \sim P_{xy}} E_{x' \sim P_x} E_{y' \sim P_y} K(x,x') \ell(y,y')$$

Thus, in term 1, $X$ and $X'$ (and $Y$ and $Y'$) should be independent, in term 2 $x, x', y, y'$ should all be independent, in term 3 $(x,x')$ and $(y,y')$ should be independent (and $(x',y')$) 
(below term 1: but not $x, y$)

—We can also test for dependence by directly computing the cross-covariance using the operator $C_{xy}$. this is called the constrained covariance:

$$COCO(P_{xy}; F, G) = \sup_{\|f\|_{\mathcal{F}}=1, \|g\|_{\mathcal{G}}=1} cov\left[f(x), g(y)\right]$$

$$= \sup \quad \mathbb{E}_{xy}\left[f(x) \otimes g(y)\right] - \mathbb{E}_x f(x) \otimes \mathbb{E}_y g(y)$$

$$= \sup \quad \mathbb{E}_{xy}\left[\langle f, \phi(x)\rangle_{\mathcal{F}} \otimes \langle g, \psi(y)\rangle_{\mathcal{G}}\right] - \mathbb{E}_x \langle f, \phi(x)\rangle \otimes \mathbb{E}$$

$$= \sup \left(f \otimes g, \mathbb{E}_{xy} \phi(x) \otimes \psi(y)\right)_{\mathcal{F} \times \mathcal{G}} - \langle f, \mu_x\rangle \langle g, \mu_y\rangle$$

$$= \sup \quad \cancel{\langle f \otimes g, \tilde{C}_{xy}\rangle_{\mathcal{F} \times \mathcal{G}}} - \langle f \otimes g, \mu_x \otimes$$

$$= \sup \left(f \otimes g, \tilde{C}_{xy} - \mu_x \otimes \mu_y\right)_{\mathcal{F} \times \mathcal{G}}$$

$$= \sup_{\|f\|_{\mathcal{F}}=1, \|g\|_{\mathcal{G}}=1} \left\langle f, C_{xy} g\right\rangle$$

which we estimate empirically using $\hat{C}_{xy}, \hat{\mu}_x, \hat{\mu}_y$. Noting again that computing $\hat{\mu}_x$, $\hat{\mu}_y$ requires dot products between $f$ and $\{x_i\}$ and $g$ and $\{y_i\}$, any components of $f$ and $g$ orthogonal to the "data space" disappear and are thus irrelevant. We can therefore express $f$ and $g$ as follows, without loss of generality:

$$f = \sum_i \alpha_i \phi(x_i) = X^{\top} \alpha \qquad g = \sum_i \beta_i \psi(y_i) = Y^{\top} \beta$$

where $\bar{\phi}(x_i) = \phi(x_i) - \frac{1}{N}\sum_i \phi(x_0)$ and equivalently for $\bar{\psi}(y_i)$.

We then have the following empirical estimate of COCO:

$$\langle X\alpha, \frac{1}{N} X H Y^T Y \beta \rangle = \alpha^T H X^T X H Y^T Y H \beta$$

since $HH = H$

$$\cancel{= \alpha^T H X^T X H H Y^T Y H \beta}$$

$$= \alpha^T H X^T X H H Y^T Y H \beta$$

$$= \alpha^T H K H H L H \beta = \frac{1}{N}\alpha^T \tilde{K} \tilde{L} \beta$$

We then solve for $\alpha, \beta$ by maximizing the following Lagrangian:

$$\mathcal{L}(\alpha, \beta, \lambda, \gamma) = \frac{1}{N}\alpha^T \tilde{K} \tilde{L} \beta - \frac{\lambda}{2}(\alpha^T \tilde{K} \alpha - 1) - \frac{\gamma}{2}\beta^T \tilde{L} \beta -$$

$$\|\beta\| = \|X H \alpha\| = 1$$
$$= \|X H \alpha\|^2$$
$$= \cancel{\|} \alpha^T H X^T X H \alpha$$
$$= \alpha^T \tilde{K} \alpha = 1$$

Differentiating, we set:

$$\begin{cases} \dfrac{\partial \mathcal{L}}{\partial \alpha} = \dfrac{1}{N} \tilde{K} \tilde{L} \beta - \lambda \tilde{K} \alpha = 0 \\[3mm] \dfrac{\partial \mathcal{L}}{\partial \beta} = \dfrac{1}{N} \tilde{L} \tilde{K} \alpha - \gamma \tilde{L} \beta = 0 \end{cases}$$

Multiplying by $\alpha^T$ for eqn 1, $\beta^T$ for eqn 2:

$$\frac{1}{N}\alpha^T \tilde{K} \tilde{L} \beta = \lambda \alpha^T \tilde{K} \alpha \qquad \lambda \alpha^T \tilde{K} \alpha = \gamma \beta^T \tilde{L} \beta$$

$$\Rightarrow \frac{1}{N}\alpha^T \tilde{K} \tilde{L} \beta = \gamma \beta^T \tilde{L} \beta \qquad \Longrightarrow \qquad \overset{1}{\lambda} \overset{1}{\underset{\smile}{\lambda = \gamma}}$$

$$\Rightarrow \begin{bmatrix} 0 & \frac{1}{N}\tilde{K}\tilde{L} \\ \frac{1}{N}\tilde{L}\tilde{K} & 0 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \lambda \begin{bmatrix} \tilde{K} & 0 \\ 0 & \tilde{L} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \qquad \text{Generalized eigenvalue problem}$$

Note that the constraints $\|f\|_{\mathcal{F}} = 1$, $\|g\|_{\mathcal{G}} = 1$ enforce smoothness in $f$ and $g$ such that COCO is insensitive to high frequency dependencies (need high sample size to detect).

- Returning to HSIC, we now ask what an unbiased estimate would be and have biased $\text{Tr}(KHLH)$ is:

$$\text{Term 1}: \langle \bar{C}_{xy}, \tilde{C}_{xy} \rangle_{HS} = \mathbb{E}_{\substack{x,y \sim P_{xy} \\ x',y' \sim P_{xy}}} k(x,x') \ell(y,y')$$

$$\simeq \left( \sum_{i=1}^{N} \sum_{j \neq i} k(x_i, x_j) \ell(y_i, y_j) \right) \frac{1}{N(N-1)}$$

Difference b/w biased and unbiased estimate is the

$$\underbrace{\frac{1}{N^2} \sum_{i,j} K_{ij} L_{ij}}_{\text{biased}} - \underbrace{\frac{1}{N(N-1)} \sum_{j \neq i} K_{ij} L_{ij}}_{\text{unbiased}}$$

$$= \frac{1}{N^2} \sum_i K_{ii} L_{ii} + \frac{1}{N^2} \sum_{j \neq i} K_{ij} L_{ij} - \frac{1}{N(N-1)} \sum_{j \neq i} K$$

$$= \frac{1}{N} \left( \frac{1}{N} \sum_i K_{ii} L_{ii} - \frac{1}{N(N-1)} \sum_{j \neq i} K_{ij} L_{ij} \right)$$

$$\mathbb{E}[\text{bias}] = \frac{1}{N} \left( \mathbb{E}_x k(x,x) \mathbb{E}_y \ell(y,y) - \frac{N(N-1)}{N(N-1)} \mathbb{E}_{x,x'} k(x,x') \mathbb{E}_y \mathbb{E}_{y'} \right)$$

$$\sim \mathcal{O}\left(\frac{1}{N}\right)$$

Term 2 : $\langle \mu_x \otimes \mu_y, \mu_x \otimes \mu_y \rangle_{F \times G} = \underset{x \sim P_x}{\mathbb{E}} \underset{x' \sim P_x}{\mathbb{E}} k(x,x') \underset{y \sim P_y}{\mathbb{E}} \underset{y' \sim P_y}{\mathbb{E}} \ell(y,y')$

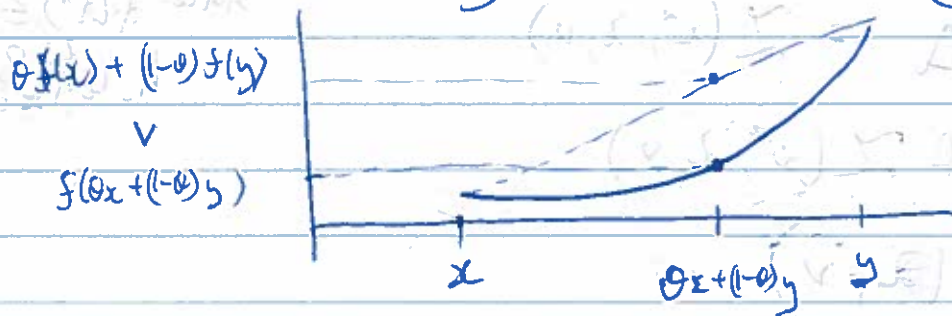$$\simeq \frac{1}{N(N-1)(N-2)(N-3)} \sum_i \sum_{j \neq i} \sum_{q \neq i,j} \sum_{r \neq i,j,q} K_{ij} L_{qr}$$

Term 3 : $\langle \tilde{C}_{xy}, \mu_x \otimes \mu_y \rangle_{F \times G} = \underset{x,y \sim P_{xy}}{\mathbb{E}} \underset{x' \sim P_x}{\mathbb{E}} k(x,x') \underset{y' \sim P_y}{\mathbb{E}} \ell(y,y')$

$$= \left[ \sum_i \sum_{j \neq i} K(x_i, x_j) \sum_{q \neq i,j} \ell(y_i, y_j) \right] \frac{1}{N(N-1)(N-}$$

# Part III: SVMs and Convex Optimization

- A set $C$ is convex if for any $x_1, x_2 \in C$, $\theta x_1 + (1-\theta) x_2 \in C$, $0 \le \theta \le 1$
- A function $f(x)$ is convex if its domain dom $f$ is a convex set and for any $0 \le \theta \le 1$,

$$f(\theta x + (1-\theta) y) \le \theta f(x) + (1-\theta) f(y)$$

$\theta f(x) + (1-\theta) f(y)$

$\vee$

$f(\theta x + (1-\theta) y)$

$x \qquad \theta x + (1-\theta) y \qquad y$

- Suppose we want to solve the following optimization problem:

$$\text{minimize} \quad f_0(x)$$
$$\text{subject to} \quad f_i(x) \le 0 \qquad i = 1, \dots, m$$
$$\qquad\qquad\qquad h_i(x) = 0 \qquad i = 1, \dots, p$$

Calling us then optimum point $f_0(x^*)$.

It turns out we can solve this by solving a different easier - convex - optimization problem, called the Lagrange dual problem:

$$\text{maximize} \quad g(\lambda, v)$$
$$\text{subject to} \quad \lambda \ge 0 \qquad\qquad (\lambda_i \ge 0 \; \forall i)$$

where $g(\lambda, v) = \inf_x \mathcal{L}(x, \lambda, v)$

$$= \inf_x \left[ f_0(x) + \sum_i \lambda_i f_i(x) + \sum_j v_j h_j(x) \right]$$

Lagrange multipliers dual variables

Lagrange dual function

$$= \sup_{\lambda \geq 0} \mathcal{L}(x^*, \lambda, \nu) \quad \left( \begin{array}{l} \text{i.e. } \lambda_i = 0 \text{ since} \\ f_i(x^*) \leq 0 \end{array} \right)$$

$$= \inf_x \sup_{\lambda \geq 0} \mathcal{L}(x, \lambda, \nu)$$

We can see this Lagrange dual problem is equivalent to our original minimization problem by noting that $g(\lambda, \nu)$ is upper bounded by $f_0(x^*)$

$$f_0(x^*) \geq f_0(x^*) + \sum_i \lambda_i f_0(x^*) + \sum_j \nu_j h_i(x^*) \qquad \text{since } \begin{array}{l} f_i(x^*) \leq \\ \lambda \geq 0 \\ h_i(x^*) = i \end{array}$$

$$= \mathcal{L}(x^*, \lambda, \nu)$$

$$\geq \inf_x \mathcal{L}(x^*, \lambda, \nu)$$

$$\geq g(\lambda, \nu)$$

We call the solution to the Lagrange dual problem $(\lambda^*, \nu^*)$ dual optimal.

We know the dual problem is convex, since $g(\lambda, \nu)$ is concave and the constraint set is convex.

$$\hookrightarrow \lambda \geq 0 \; ; \; g(\lambda, \nu) > -\infty$$

Two cases are possible
when replacing the original opt. $(\lambda, \nu)$ is dual feasible
problem w/ the dual problem:

$$g(\lambda^*, \nu^*) \leq f_0(x^*) \qquad \text{weak duality}$$

$$g(\lambda^*, \nu^*) = f_0(x^*) \qquad \text{strong duality}$$

The conditions under which strong duality holds, are called constraint qualifications

(put simply, there exists an $\bar{x}$ that satisfies all the constraints)

(sufficient, not necessary conditions for strong duality)

I think!

One example is:

- primal problem is convex: i.e. $f_i(x)$ convex and $h(x) = Ax + b$

- <u>Slater's condition holds:</u> there exists some strictly feasible point $\tilde{x}$ s.t. $f_i(\tilde{x}) < 0$ $\forall i$ and $h(\tilde{x})$:

Note that when strong duality holds (i.e. the dual problem is equal to the primal),

$$f_0(x^*) = g(z^*, v^*)$$

$$\geq \inf_x \; f_0(x) + \sum_i z_i^* f_i(x) + \sum_i v_i^* h_i(x)$$

$$\leq f_0(x^*) + \sum_i z_i^* f_i(x^*) + \sum_i v_i^* h_i(x^*)$$

$$\leq f_0(x^*) \implies \sum_i z_i^* f_i(x^*) = 0$$

This implies <u>complimentary slackness:</u>

$$z_i^* \neq 0 \implies f_i(x^*) = 0$$

$$f_i(x^*) < 0 \implies z_i^* = 0$$

Lastly, if $f_0, \{f_i\}$ and $\{h_i\}$ are all differentiable, then

$$\nabla f_0(x^*) + \sum_i z_i^* \nabla f_i(x^*) + \sum_j v_j \nabla h_i(x^*) = 0$$

Putting all the above conditions together gives the <u>KKT conditions!</u>

$f_i(x) \leq 0$      $z_i f_i(x) = 0$

$h_i(x) = 0$      $\nabla f_0(x) + \sum_i z_i \nabla f_i(x) + \sum_i v_i \nabla h_i(x) = 0$

$z_i \geq 0$

~~When strong duality holds, Slater's condition holds,~~
~~valid~~

When an optimization problem is convex and
Slater's condition holds (i.e. strong duality holds)
if $f_0, \{f_i\}, \{h_i\}$ are differentiable then the
KKT conditions are necessary and sufficient
for global optimality, i.e., a solution $x^*$ that satisfies
the KKT conditions is a global optimum.

- The <u>representer theorem</u> says that
  the solution to ~~an~~ arbitrary optimization problem
  $$f^* = \underset{f \in H}{\text{argmin}} \; L_y(f(x_1), \dots, f(x_N)) + \Omega\left(\|f\|_H^2\right)$$
  where $L_y$ is parametrized by $\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ and $\Omega(\cdot)$ is
  non-decreasing, takes the form:
  $$f^* = \sum_{i=1}^{N} \alpha_i K(x_i, \cdot)$$
  where $K$ is the kernel corresponding to the RKHS $H$
  where $f$ lives.

Pf. Let $f = f_s + f_\perp$, where $f_s$ ~~is the projection~~
  is the projection of $f$ onto the ~~subspace~~ space spanned
  by $\{K(x_i, \cdot)\}_{i=1}^{N}$ and $f_\perp$ is the orthogonal error.
  By reproducing property of RKHS $H$,
  $$f(x_i) = \langle f, x_i \rangle_H = \langle f_s, x_i \rangle_H + \langle f_\perp, x_i \rangle_H$$
  $$= \langle f_s, x_i \rangle_H$$

So $L_y(f(x_1), \dots, f(x_N)) = L_y(f_s(x_1), \dots, f_s(x_N))$

$$\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \langle f_s, f_s \rangle_{\mathcal{H}} + \langle f_\perp, f_\perp \rangle_{\mathcal{H}} - 2\langle f_s, f_\perp \rangle_{\mathcal{H}}$$

$$= \|f_s\|_{\mathcal{H}}^2 + \|f_\perp\|_{\mathcal{H}}$$

If $\Omega(\cdot)$ is strictly increasing, then its minimum is achieved when $\|f_\perp\|_{\mathcal{H}}^2 = 0$ and $\|f_s\|_{\mathcal{H}}^2$ is minimized, leaving the optimization problem

$$f^* = \underset{f_s = \sum_{i=1}^{N} \alpha_i k(x_i, \cdot)}{\arg\min} L_y\left(f_s(x_1), \ldots, f_s(x_N)\right) + \Omega\left(\|f_s\|_{\mathcal{H}}^2\right)$$

– In **support vector classification** we want to find a hyperplane that separates data of two different classes within the data space itself.



The best such hyperplane is the one that maximizes the distance b/w the margins while enforcing perfect classification.

Let $w$ = vector perpendicular to hyperplane. We want:

① $w^T x_i + b \geq 1 \quad \forall i : y_i = 1$

$w^T x_i + b \leq -1 \quad \forall i : y_i = -1$

$\Rightarrow y_i(w^T x_i + b) \geq 1$

② For $x^+, x^-$ on opposite margins,

$$\text{maximize} \quad \frac{x^{+T} w}{\|w\|} - \frac{x^{-T} w}{\|w\|} = \frac{1}{\|w\|} - \frac{-1}{\|w\|} = \frac{2}{\|w\|}$$

So, we want to solve

$$\text{maximize} \frac{2}{\|w\|} \Big/ \text{minimize} \frac{1}{2}\|w\|^2$$

$$\text{subject to} \quad y_i(w^T x_i + b) \geq 1 \quad \forall i$$

However, there will rarely exist a hyperplane that exactly separates the data (i.e. impossible to get $y_i(w^T x_i + b) \geq$ for all $x_i$), so we soften the constraints with an error term $\xi_i$, which we minimize as well:

$$\min_{w,b} \left( \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{N} \xi_i \right) \qquad \left( \begin{array}{l} \text{controls trade-off} \\ \text{b/w flexibility and} \\ \text{accuracy} \end{array} \right)$$

$$\text{subject to} \quad y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

Giving us the following Lagrangian dual function:

~~$$\mathcal{L}(w,b,\xi,\lambda,v)$$~~

$$\mathcal{L}(w,b,\xi,\lambda,v) = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N} \xi_i + \sum_{i=1}^{N} \lambda_i (1 - y_i(w^T x_i + b) - $$

$$- \sum_{i=1}^{N} v_i \xi_i$$

· Noting that each of our constraints are convex and there always exists a set of $\{w, b, \{\xi_i\}\}$ that satisfies them (i.e. Slater's condition holds), since the objective and constraint functions are differentiable we need only solve for the KKT conditions to find the global optimum

prove that strong duality holds and we can solve for the KKT conditions and optimize the dual

$$\lambda_i \geq 0, \quad v_i \geq 0$$

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^{N} \lambda_i y_i x_i = 0 \iff w = \sum_{i=1}^{N} \lambda_i y_i x_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = -\sum_i \lambda_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \lambda_i - v_i = 0$$

$$\iff \lambda_i = C - v_i$$

Since $z_i, v_i \geq 0$, we have $0 \leq z_i \leq C$.
Then, by complimentary slackness, we can work out the following three possible cases:

For $z_i = C$,
$$v_i = 0 \Rightarrow \xi_i \geq 0 \qquad \text{i.e. } x_i \text{ lies inside}$$
$$y_i(w^T x_i + b) = 1 - \xi_i \leq 1 \qquad \text{the margin}$$

For $0 < z_i < C$,
$$v_i \geq 0 \Rightarrow \xi_i = 0 \qquad \text{i.e. } x_i \text{ lies on}$$
$$y_i(w^T x_i + b) = 1 \qquad \text{the margin}$$

For $z_i = 0$,
$$v_i = C \Rightarrow \xi_i = 0 \qquad \text{i.e. } x_i \text{ correctly outside}$$
$$y_i(w^T x_i + b) \geq 1 \qquad \text{the margins}$$

In sum, ~~and solution is~~ only points on or inside the margin with $z_i > 0$ contribute to ~~the support vectors~~ $w = \sum z_i y_i x_i$, their contribution bounded by $C$. These points are thus called the support vectors. Us note also that this is a sparse solution, since most $x_i$ will be outside the margin so $z_i = 0$.

Given that strong duality holds, we can find $z_i$'s by ~~appropriately~~ maximizing the Lagrangian dual function, ~~plugging in our KKT~~ plugging in our KKT conditions to simplify.

$$g(\alpha, \nu) = \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j + C \sum_{i=1}^{N} \xi_i$$

$$+ \sum_i \lambda_i \left( 1 - y_i \left( \sum_j \lambda_j y_j x_j^T x_i + b \right) - \xi_i \right)$$

$$- \sum_i \nu_i \xi_i$$

$$= \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j + C \sum_i \xi_i$$

$$- \sum_i \lambda_i \xi_i + \sum_i \lambda_i \quad - \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j$$

$$- b \sum_i \lambda_i y_i \quad - \sum_i \nu_i \xi_i$$

$$= -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_i \overbrace{(C - \lambda_i)}^{\nu_i} \xi_i - \sum_i \nu_i \xi_i + \sum_i$$

$$= -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_i \lambda_i = g(\lambda)$$

Thus, to find the support vector we solve

$$\text{maximize} \quad g(\lambda) = -\frac{1}{2} \|w\|^2 + \sum_i \lambda_i$$

$$\text{subject to} \quad 0 \leq \lambda_i \leq C$$

We can the estimate $b$ by ~~solving~~ solving $y_i(w^T x_i + b) = 1$
for our $x_i$ on the margin (or ~~by~~ averaging over all $x_i$ on the m

— Since $C$ is hard to interpret, we can reparametrize with a new parameter $\nu$. This is called $\nu$-SVM:

$$\min_{w, \rho, \xi} \frac{1}{2}\|w\|^2 - \nu\rho + \frac{1}{N}\sum_{i=1}^{N}\xi_i$$

$$\text{subject to } y_i(w^Tx_i + b) \geq \rho - \xi_i$$

$$\xi_i, \rho \geq 0$$

Again Slater's condition holds, and our constraints are convex, so we go on to write out the KKT conditions and then optimize the Lagrangian!

$$\lambda_i \geq 0 \qquad \xi_i \geq 0 \qquad \frac{\partial \mathcal{L}}{\partial w} = w - \sum \lambda_i y_i x_i = 0$$

$$\alpha_i \geq 0 \quad y_i(w^Tx_i + b) \geq \rho - \xi_i \qquad \Rightarrow w = \sum_i \lambda_i y_i x_i$$

$$\gamma \geq 0 \qquad \rho \geq 0 \qquad \frac{\partial \mathcal{L}}{\partial b} = -\sum_i \lambda_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = \frac{1}{N} - \lambda_i - \alpha_i = 0$$

$$\Rightarrow \lambda_i + \alpha_i = \frac{1}{N}$$

$$\mathcal{L}(w, \rho, \xi, \lambda, \alpha, \gamma)$$

$$= \frac{1}{2}\|w\|^2 - \nu\rho + \frac{1}{N}\sum_i \xi_i + \sum_i \lambda_i\left(\rho - y_i(w^Tx_i + b) - \xi_i\right)$$

$$- \sum_i \alpha_i \xi_i \qquad - \gamma\rho$$

$$\frac{\partial \mathcal{L}}{\partial \rho} = -\nu + \sum_i \lambda_i - \gamma = 0$$

$$\Rightarrow \nu = \sum_i \lambda_i - \gamma$$

$g(\alpha, \mu, \gamma)$

$$\cancel{g(\alpha,\mu,\gamma)} = -\frac{1}{2}\sum_{i,j}z_i z_j y_i y_j x_i^T x_j - \left(\sum_i z_i\right)\rho + \cancel{\rho}\rho$$

$$+ \frac{1}{N}\sum_i \xi_i + \sum_i z_i \rho - \sum_i z_i \xi_i$$

$$+ \cancel{\rho\rho} - \sum_i \alpha_i \xi_i - \gamma\rho$$

$$\boxed{= -\frac{1}{2}\sum_{i,j} z_i z_j y_i y_j x_i^T x_j = g(z)}$$

So, ~~writing the Lagrange dual problem~~

~~minimize $g(z) = -\frac{1}{2}\sum_{i,j}z_i z_j y_i y_j x_i^T x_j$~~
~~subject to~~

~~By complementary slackness~~

So, Lagrange dual problem becomes:

minimize $g(z)$

subject to $0 \leq z \leq \frac{1}{N}$

How do we interpret $\gamma$?

- assume $\rho \geq 0 \implies \gamma = 0 \implies \cancel{\rho}\gamma = \sum_i z_i$
- by complementary slackness,

  If $\xi_i > 0 \implies \alpha_i = 0 \implies z_i = \frac{1}{N}$

  If $\xi_i = 0 \implies \alpha_i \geq 0 \implies z_i \leq \frac{1}{N}$

$\implies$ For $N(z) = \{z_i = \frac{1}{N}\}$, $\sum_{i \in N(z)} z_i + \sum_{j \in M(z)} z_j < \frac{|N(z)| + |M(z)|}{N} \leq \sum_i z_i = \gamma$

$M(z) = \{0 < z < \frac{1}{N}\}$

- We can "Kernelize" SVM by moving into a feature space $\mathcal{H}$ : $w = \sum_i c_i K(x_i, \cdot)$, $w \in \mathcal{H}$. Writing down the objective as follows:

$$w^* = \underset{w \in \mathcal{H}}{\arg\min} \; \frac{1}{2}\|w\|_{\mathcal{H}}^2 + C \sum_i \xi_i$$

$$= \underset{w \in \mathcal{H}}{\arg\min} \; \frac{1}{2}\|w\|_{\mathcal{H}}^2 + C \sum_i \left[ \max\left(0, 1 - y_i \left(\langle w, K(x_i, \cdot)\rangle_{\mathcal{H}} + b\right)\right)\right]$$

rectification

$$= \underset{w \in \mathcal{H}}{\arg\min} \; \Omega(\|w\|_{\mathcal{H}}) + F_y\left(\langle w, K(x_1, \cdot)\rangle, \dots, \langle w, K(x_N, \cdot)\rangle\right)$$

we recognize that $\Omega(\cdot) = \frac{1}{2}(\cdot)^2$ is non decreasing so the representer theorem applies, giving us

$$w^* = X\beta$$

Thus the minimization problem becomes

$$\text{minimize} \quad \frac{1}{2}\beta^T K \beta + C \sum_i \xi_i$$

$$\text{subject to} \quad \xi_i \geq 0$$

$$y_i\left(\sum_j \beta_j K(x_i, x_j) + b\right) \geq 1 - \xi_i$$

Since $K$ is positive definite, the objective function and constraints are convex so Slater's condition gives us strong duality.

We can thus instead optimize the dual, as before:

$$\text{minimize } g(\lambda) = \frac{1}{2}\sum_{i,j} \lambda_i \lambda_j \, y_i y_j \, K(x_i, x_j) + \sum_i \lambda_i$$

$$\text{subject to } 0 \leq \lambda_i \leq C$$

— My convex optimization recipe:

1. write down objective and constraints
2. Check constraints don't conflict
   if not → Slater's condition holds
3. Check if constraints and objective are convex
   if so → strong duality holds
4. Write down Lagrangian
5. write down KKT conditions ~~of duality~~
   (need to differentiate Lagrangian)
6. Expand Lagrangian and simplify using
   KKT conditions, plugging in where possible
7. Rewrite optimization problem in terms
   of the Lagrangian dual subject to
   constraints implied by KKT conditions
8. Use complimentary slackness conditions
   to interpret different parameter values