

الدكتور انهارد هيتكلم علي ال kernel methods .. بتسمحك انك ت scale شوية حاجات لما بتشتغل مع linear models فاحنا لما كنا بنتكلم علي linear models و كان في طريقه عشان نشغل علي non-linear models تعال نشغل علي phi mapping تودينا من ال original space to a new space وهناك اعمل ال linear algebra بتاعتك ...

تعال نفكر شوية حاجات ..

Generalized linear models:

- Work with fixed non-linear basis functions
- Hypothesis space is limited
- Optimization is easier (usually convex)

Neural Networks:

- Adaptive non-linear basis functions
- Richer hypothesis space
- Harder to optimize (Draw back) (Usually optimization is non-convex)

ده كان ملخص كدا لل 2 important paradigms في الماشين ليرننج ..

طيب في ال generalized linear models احنا بنحتاج نجيب fixed basis functions .. المشكله في ال cost omputation .. في اسمها ال dual ... وديه بتقول ان ال complexity of the computation doesn't have to depend on #of basis functions we have but depends on amount of data we have

فلو كان عندنا low data .. بيقا التكنيكس ديه ممكن تبقا فعّاله عشان انت هتقدر ت map the data .. والمابينج هيعتمد علي حجم الداتا .. وده بيتم بال dual trick في معظم الألوورزمز لما بنروح ل new space هي دايماً من نوع $\phi(x)^T * \phi(x')$ اللي هو ان ال computation بيبيقا dot product between ... فلو هو دايماً dot product فانت ممكن ت accelerate الموضوع ده ... قال kernel function هتبقا هي ال pair of points انت مش محتاج تعرف ايه x and x prime في ال new space اللي عندنا .. فلو انت عارف الأوتوتوت بتاع ال kernel function لكل pair of points انت مش محتاج تعرف ايه ال underlying feature space defined by $\phi(x)$.. اللي هي يعني $k(x, x')$ وبعدها هنبيين ان ممكن نعرّف kernel functions اللي بتبقا سريعه جداً عشان ت evaluate و مش محتاجه explicit computation with $\phi(x)$

كمثال تعال نشوف ال linear regression و تعال نشوف ايه اللي بيحصل (in terms of computation) ... فالأول هنشوف ان كل ال computation بتبقا $\phi^T \phi$ وده ال kernel .. وبعدين هنشوف ان ممكن نعرّف ال kernel بحيث ان ال computation is independent of the dimensionality of the space فبالنسبه لل linear regression ده كان الهدف بتاعنا .. احنا بنقل ال squared loss و بنضيف ترم ال regularization ... فاحنا خدنا المشتقه وساوينها بالصفر .. وفي الآخر بدل ما تعزل ال w وتجب الحل .. اكتب الحل بطريقه مختلفه ...

Dual Representations

- Recall linear regression objective

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [\mathbf{w}^T \phi(\mathbf{x}_n) - y_n]^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Solution: set gradient to 0

$$\nabla E(\mathbf{w}) = \sum_n (\mathbf{w}^T \phi(\mathbf{x}_n) - y_n) \phi(\mathbf{x}_n) + \lambda \mathbf{w} = 0$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_n (\mathbf{w}^T \phi(\mathbf{x}_n) - y_n) \phi(\mathbf{x}_n)$$

∴ **w is a linear combination of inputs in feature space**

$$\{\phi(\mathbf{x}_n) | 1 \leq n \leq N\}$$

لما هنتجعي نعوض عن الداليليو ب الإكسبرشن اللي قدامك ده .. و هتقول ان الغاي هي الماتركس بتاعت كل النقط اللي عندك ... فالغاي بتاعت إكس 1 .. هي ال 1st data point
فالفاي ماتركس ديه كلها هي ال concatenation of data points
اللي تقدر تعملو ناو اتنا تعامل ال a على اساس انو variable بدل W .. عشاش الفاي هي الداتا في النيو سيمس .. ديه حاجه عندنا اصلا .. فلما هت vary هو فعلياً كزناك بت vary
a .. فانت كمان ممكن تحط في اعتبارك ال a كزناك new variable .. فال dual trick هو انك تعمل ال optimization into a different space فاللي بيحصل ان عندنا داتا في
السيبس و عندنا ماينج فاي بتودينا للسيبس الثاني .. بس كمان انت عندك فاريل W هتعملو mapping ل dual space فال a هنا هي مجرد w equivalent بس في سيمس
مختلف

- Substitute $\mathbf{w} = \Phi \mathbf{a}$
- Where $\Phi = [\phi(x_1) \phi(x_2) \dots \phi(x_N)]$

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} \quad \text{and} \quad a_n = -\frac{1}{\lambda} (\mathbf{w}^T \phi(x_n) - y_n)$$

- Dual objective: minimize E with respect to \mathbf{a}

$$E(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \underbrace{\Phi^T \Phi}_{\mathbf{K}} \Phi^T \Phi \mathbf{a} - \mathbf{a}^T \Phi^T \Phi \mathbf{y} + \frac{\mathbf{y}^T \mathbf{y}}{2} + \frac{\lambda}{2} \mathbf{a}^T \Phi^T \Phi \mathbf{a}$$

الدكتور بيقولك ان لما بتعمل ال computation to the dual space احنا عندنا different complexit for the computation في المسألة الأصلية اللي عندنا
بالنسبة لل w كانت بتعتمد على عدد ال basis functions .. فلما بتعمل ال optimization with respect to w اني عندني variable W for every basis function .. بس لما هتشتغل في ال
dual space .. هتشتغل على ال a بدل ال w .. فالكوميلاستي هتتعتمد على كمية الداتا مش ال dimensionality of space ...

- Prediction: $y_* = \phi(x_*)^T \Phi \mathbf{a}$
 $= k(x_*, X)(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$
- Linear regression where we find dual solution \mathbf{a} instead of primal solution \mathbf{w} .

Complexity:

- Primal solution: depends on # of basis functions
- Dual solution: depends on amount of data
 - Advantage: can use very large # of basis functions
 - Just need to know kernel k

ازاي نوصل لل complexit بتعتمد على حجم الداتا هيعتمد على طريقه نحسب بينها ال kernels اللي هم ال gram matrix من غير ما نضطر اننا ن
refer to the points into the new space .. فانت محتاج بيكا عندك kernel K .. ديه بتسمحك انك تلاقى ال dot products الي انت عاوزهم دول
من غير ما تحسب ال dot product ... في طرق كتيره تقدر بينها تعرف kernel functions اللي هت implicitly correspond to dot products

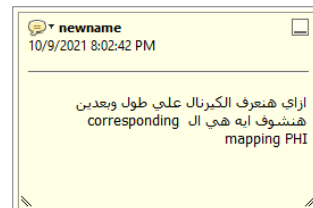
..... انت مش محتاج تعمل ال dot products هو في طريقه مباشره اللي هي ال kernel ... فانت عاوز ال Gram matrix K وانت عارف انها الفاي ترانسبوز فاي .. فبدل ما تروح تحسب الدوت برودكت .. علي طول K specify .. فانا لو عملت كدا هنا في شوية فانكشنز تمام هيقدر و يعملو كدا .. بس في فانكشنز مش هينفع .. فهي مش هتبقا arbitrary function .. لازم تختار $\phi^T \phi$ functions that have an implicit correspondence with $\phi^T \phi$..

..... فالفانكشنز اللي ليها implicit correspondence هيبقا لازم يكونوا positive semi-definite ... الماتركس بتبقا positive semi-definite لما تكون تقدر ت factor الماتركس ديه ل ϕ one matrix ϕ و تضربها في نفسها هيديلك ال K ثاني .. فعشان كدا ال Gram matrix or kernel function بتحتاج تبقا positive semi definite .. معناها بس ان لازم بيقا في ماتركس فاي تقدر تضربها ف نفسها .. فلو بنتكلم linear algebra هي هي انك تقول ان ليها eigenvalues أكبر من او بتساوي الصفر ..

Constructing Kernels

- Two possibilities:
 - Find mapping ϕ to feature space and let $K = \phi^T \phi$
 - Directly specify K
- Can any function that takes two arguments serve as a kernel?
- No, a valid kernel must be positive semi-definite
 - In other words, k must factor into the product of a transposed matrix by itself (e.g., $K = \phi^T \phi$)
 - Or, all eigenvalues must be greater than or equal to 0.

Example



- Let $k(x, z) = (x^T z)^2$

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

$$k(x, z) = (x^T z)^2$$

$$= (x_1 z_1 + x_2 z_2)^2$$

$$= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2$$

$$= (x_1^2, 2x_1 x_2, x_2^2)$$

$$\begin{bmatrix} z_1^2 \\ 2z_1 z_2 \\ z_2^2 \end{bmatrix} = \phi(x)^T \cdot \phi(z)$$

ϕ_1 ϕ_2 ϕ_3

طبيب في فكره عشان نقدر نبني ال kernels هو ان في basic rules نقدر نمشي وراها .. اللي هو لو عندي kernels اصلا من الأول .. فاقدر أ compose أكثر من كيرنال عشان يبقا عندي kernels ثانيين

Rules to construct Kernels

- Let $k_1(x, x')$ and $k_2(x, x')$ be valid kernels
 - The following kernels are also valid:
 - $k(x, x') = ck_1(x, x') \quad \forall c > 0$
 - $k(x, x') = f(x)k_1(x, x')f(x') \quad \forall f$
 - $k(x, x') = q(k_1(x, x'))$ q is polynomial with coeffs ≥ 0
 - $k(x, x') = \exp(k_1(x, x'))$
 - $k(x, x') = k_1(x, x') + k_2(x, x')$
 - $k(x, x') = k_1(x, x')k_2(x, x')$
 - $k(x, x') = k_3(\phi(x), \phi(x'))$
 - $k(x, x') = x^T A x' \quad A$ is symmetric positive semi-definite
 - $k(x, x') = k_a(x_a, x'_a) + k_b(x_b, x'_b)$
 - $k(x, x') = k_a(x_a, x'_a)k_b(x_b, x'_b)$
- where $x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$

في شوية كيرنالز موجودين في ارض الواقع من ضمنهم ال polynomial kernel ... هنا انت بتاخذ الدوت بروضكت في الأوريجنال سبيس وبعدين ترفعها لباور .. لما ترفع لباور M .. هتجيب polynomial expression of degree M .. الفيتشر سبيس هو ال all degree M products of entries in x ... لو عندك ال x and x prime are 2 images .. في الحالة ديه ال feature space هيبقا كل الضرب بتاع ال M pixel intensities تعال نرجع للمثال اللي فات ... كان في باور 2 .. فهنا ال $M = 2$.. و ال implicit feature space هيبقا monomials of degree M

الكيرنالز هي a special construction that allows dot product between pairs of data points ... وده مفيد لو هنتشتغل في ال generalized linear models و هنا انت بتبدأ ب space وبعدين تعمل mapping phi لتوديك ل new space وبعدين توديك تعمل regression and classification هناك

المشكلة اللي شفناها ان لما بتروح لل new space .. بيبقا ليه high dimensionality وبتبقا عاوزلو lots of basis function هتفتدفع التمن computation cost ... انما بالكيرنالز .. كل الحسابات عباره عن dot product في ال new space .. فال kernel بتبقا في ال new space ايه هي ال dot product ما بين pairs of points و دلوقت لما بتقدر تحدد ده علي طول من غير ما تضطر تروح لل new space و تعمل هناك ال dot product فانت تقدر ت save the computation و تدفع التمن بتاع ال dimensionality في ال original space ...

CS480/680 Lecture 11: Kernel Methods

$$k(x, x') = (x^T x')^2$$

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad x' = \begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix}$$

$$= (x_1 x'_1 + x_2 x'_2)^2$$

$$= x_1^2 x'^2_1 + 2x_1 x'_1 x_2 x'_2 + x_2^2 x'^2_2$$

$$= \begin{pmatrix} x_1^2 & \sqrt{2} x_1 x_2 & x_2^2 \end{pmatrix} \begin{pmatrix} x'^2_1 \\ \sqrt{2} x'_1 x'_2 \\ x'^2_2 \end{pmatrix}$$

1 prime X 2 prime and then X 2 prime square

تعال نتكلم علي مثال لو خدنا ال x and x' بيقو صورتين .. الإكس هبيقا فيكتور جواه pixel intensities .. لما هتروح لل new feature space عن طريق انك تاخذ ال x transpose x' وترفعهم لباور .. انت هتجيب فيتشرز تانيه .. هتبقا الضرب بتاع every combination of M pixels together where some pixels might repeat .. فالجمال هنا ان دلوقت انت بقا عندك richer space وبعدين تتمني ان في ال richer space عندك ده تقدر تلاقي فانكشن اللي هي linear و implicitly nonlinear في الأوريجنال سبيس طيب ده بيشتغل كويس لو انت عاوز فيتشرز ليها degree M .. طب بالنسبة للفيتشرز بتاعت كل الدرجيز up to M .. فهنا انت بتضيف ثابت اسمو C .. فده كذا انت بتاخذ في الاعتبار عندك كل ال degrees up to M ...

Common Kernels

- Polynomial kernel: $k(x, x') = (x^T x')^M$
 - M is the degree
 - Feature space: all degree M products of entries in x
 - Example: Let x and x' be two images, then feature space could be all products of M pixel intensities
- More general polynomial kernel:

$$k(x, x') = (x^T x' + c)^M \text{ with } c > 0$$
 - Feature space: all products of up to M entries in x

مثال علي كذا:

$$\begin{aligned}
 k(x, x') &= (x^T x' + c)^2 \\
 &= (x_1 x'_1 + x_2 x'_2 + c)^2 \\
 &= x_1^2 x_1'^2 + 2x_1 x'_1 x_2 x'_2 + x_2^2 x_2'^2 + 2x_1 x'_1 c + 2x_2 x'_2 c + c^2 \\
 &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2, \sqrt{2} c x_1, \sqrt{2} c x_2, c) \\
 &\quad (x_1'^2, \sqrt{2} x'_1 x'_2, x_2'^2, \sqrt{2} c x'_1, \sqrt{2} c x'_2, c)^T
 \end{aligned}$$

هنشوف مثال علي كيرنال تاني .. جاوسين كيرنال

Common Kernels

- Gaussian Kernel: $k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$

- Valid Kernel because:

- Implicit feature space is infinite!

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

$$= e^{-\underbrace{x^T x}_{\text{red}}/2\sigma^2} e^{\underbrace{x^T x'}_{\text{red}}/\sigma^2} e^{-\underbrace{x'^T x'}_{\text{red}}/2\sigma^2}$$

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

$$= e^{-x^T x/2\sigma^2} e^{x^T x'/\sigma^2} e^{-x'^T x'/2\sigma^2}$$

$x^T x'$ is a valid kernel **by rule** 8 when $A = I$

$x^T x'$	1	1	1	1	1
$x^T x/\sigma^2$	1	1	1	1	1
$x'^T x'/\sigma^2$	1	1	1	1	1
$e^{-x^T x/2\sigma^2}$	1	1	1	1	1
$k(x, x')$	1	1	1	1	1