

أول ماشين ليرننج تكنيك .. وهو ال Nearest neighbor ... مراجعه سريعه علي الي خذناه المره اللي فاتت ... عندك شوية إنبوت X و عندك f function احنا مش عارفينها .. والسؤال هل احنا نقدر نشوف فانشكن f تقدر تمااب الانبوت للاوتبوت .. فاحنا بنسيرش علي h hypothesis عشان نقدر نحدد ال f ... سؤال سريع ... ده اسمو inductive learning ... ايه الفرق ما بينو و ما بين ال deductive learning .. محدش عارف الاجابه .. ال ديداكثيف هو ان عندك شوية رولز وشوية كونسبتس .. و بيبقا عندك شوية أمثله وبعدين يقولك حط الروولز علي الأمثله و ت deduce ايه اللي يتعمل في الأمثله اللي عندك .. الإنداكثيف ليرننج هو العكس .. انت عندك شوية أمثله و انت عاوز تعرف ايه هي ال underlying functions اللي بت exaplin الأمثله ديه ... لما هنتكلم علي ال supervised learning في 2 نوعين من المشاكل .. اول حاجه هو الكلاسيفيكيشن و ثاني حاجه هو ال regression ..

Supervised Learning

- Two types of problems
 1. Classification
 2. Regression
- NB: The nature (categorical or continuous) of the domain (input space) of f does not matter

في الكلاسيفيكيشن: ال range (output space) of function f is categorical

في ال ريجريشن: ال range (Output space) of function f is continuous

ركز ان النوع هنا بيفرق علي الاوتبوت .. مالوش دعوه بالانبوت ..

تعال نبص علي أمثله .. هنبص علي setting .. عندك outdoor sport .. وانت عندك implicit function .. وبناءً علي الجو هتقول إنك يا تستمتع بيه او مش هتستمتع بيه ..

Classification Example

- Problem: Will you enjoy an outdoor sport based on the weather?
- Training set:

Sky	Humidity	Wind	Water	Forecast	EnjoySport
Sunny	Normal	Strong	Warm	Same	yes
Sunny	High	Strong	Warm	Same	yes
Sunny	High	Strong	Warm	Change	no
Sunny	High	Strong	Cool	Change	yes

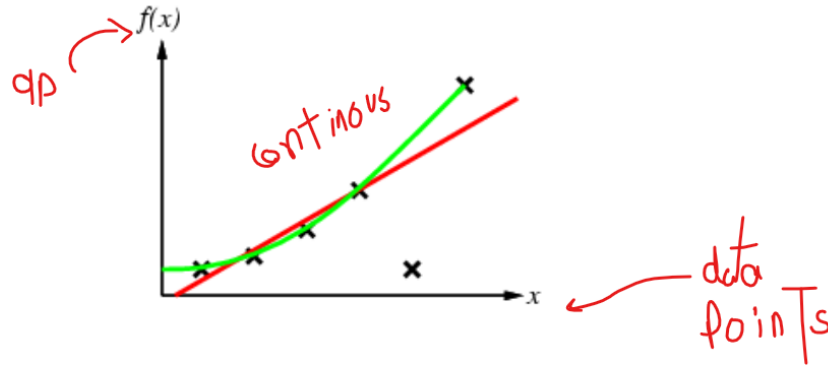
$\underbrace{\text{Sky, Humidity, Wind, Water, Forecast}}_x \quad \downarrow \quad f(x)$

- Possible Hypotheses:
 - $h_1: S = \text{sunny} \rightarrow \text{enjoySport} = \text{yes}$
 - $h_2: W = \text{cool or } F = \text{same} \rightarrow \text{enjoySport} = \text{yes}$

الفانكشن f هي مش عندك .. انت بتعملها induce ..

Regression Example

- Find function h that fits f at instances x



تعال ناخذ شوية أمثله من علي أرض الواقع ..

اول مثال عندنا هو ال **spam detection** .. لادكتور سال سؤال .. ايه هو الانبوت بتاع ال spam detector .. حد قالو اغلب الانبوت بيقا تكست ..
 سال ثاني ازاي نقدر نربرزنت التكست .. حد قالو هنعمل فيكتور و عندنا أكثر من representation .. حاجه مثلاً زي one hot vector او vector
 .. embedding

قال Domain هيقا text ... اللي هو vector of words

ال range هيقا {spam, legit}

ثاني مثال عندنا هو ال **Stock prediction**

ايه هو الانبوت: Time series of prices

ايه هو الاوتبوت space: {up, down} or \mathbb{R}^+ (all positive numbers)

خد بالك هو لو regression problem انت هيقا عندك ارقام مش بس ال Real .. بس في اغلب الوقت ممكن تعمل كذا ..

ثالث مثال هو ال **speech recognition** ..

ايه هو الانبوت: Voice (time series of audio signal)

ايه هو الاوتبوت: هي ممكن تبقا هو مين اللي بيتكلم .. او انك تبقا مهتم بايه هو اللي بيتقال .. لو انت خدت في الاعتبار الكلمات ... هيقا ايه هو نوع المشكله

ايه هو النوع: Classification عشان الحاجات اللي بتقولها هي discrete ... بس لو حد كان اتعرض لل NLP .. الكلمات بتبقا encoded using some
 embedding ... هل كذا بقت regression ولا لا؟ لا .. لأن في الآخر هي حاجه واحده اللي هتطلع مش حاجات كتير ورا بعضها ..

رابع مثال هو ال **digit recognition**:

الانبوت في الحاله ديه هي images

الاوتبوت هو الارقام بتاعتنا

والنوع هو classification

خامس مثال housing evaluation :

الانبيوت هو معظم ال features بتاعت البيوت يعني .. House features

الانبيوت هو ال positive real number R^+

ديه هتبقا regression problem

آخر واحد هيبقا ال Weather prediction

الانبيوت هو sensor data من ضمنها ال .. satellite imagery

الانبيوت هو درجة الحرارة الي هي رقم حقيقي و {Rain, sun}

نوعها : ريجريشن لدرجة الحرارة و classification للشمس و المطر

يلا نكمل ... الدكتور بيقول ان احنا عشان ندور علي ال Hypothesis هحتاج ندور في space .. اسمو ال hypothesis space

Hypothesis Space

- Hypothesis space H
 - Set of all hypotheses h that the learner may consider
 - Learning is a search through hypothesis space
- Objective: find h that minimizes
 - Misclassification
 - Or more generally some error function with respect to the training examples
- But what about unseen examples?

في ال space ده انت ممكن تدور علي optimization .. وهو ده ال objective .. احنا لما بن minimize h .. في ال unseen هم جايين من نفس ال distribution .. الهدف في الآخر تلاقي hypothesis تقدر ت generalize ..

Generalization

- A good hypothesis will generalize well
 - i.e., predict unseen examples correctly
- Usually ...
 - Any hypothesis h found to approximate the target function f well over a **sufficiently large set of training examples** will also approximate the target function well over any unobserved examples

في حاجه مهمه .. لما تلاقي hypothesis .. اتأكد انها تتوافق مع كل ال examples .. ديه حاجه انت عاوزها .. انت عاوز تلاقي hypothesis ماسكه نفسها معظم الوقت .. بس مش دايماً بتبقا متاحه .. ده هتلاقيه فعلياً مش متاح في حالات اللي بيبقا فيها insufficient hypothesis space .. فمثلاً يعني لو عندي function f والمعادله بتاعتها زي منت شايف في السلايد .. ديه فانكشن انت بتحاول تكتشف هي ايه .. و حالياً قاعد يايعيني بتدور علي الفانكشن في ال space .. أين أنتي يا فانكشن .. بس اللي انت مش واخد بالك منو هو ان السبيس اللي بتدور فيه هو space of polynomials of finite degree .. فانت مش هتعرف تلاقي اللي انت عاوزو في ال space ده .. الدكتور سأل سؤال .. حد يعرف ليه الفانكشن اللي هي $ax + b + x\sin(x)$ ليه مش هتلاقها في ال space of polynomials of finite degrees?

المشكله في ال sin function .. لو عملتلها representation by taylor series of expansion .. هتلاقي ان عندك series طويله .. بس ال degrees بتاعت ال polynomials اللي هتطلع من ال expansion هتكون و هتكون منك infite large degree for the polynomial .. فانت مش هتقدر ت represent $\sin(x)$ function ..

تاني مشكله هتظهر هو ال Noisy data ... مش هتلاقي علي طول كذا general consistent hypothesis

Inductive Learning

- Goal: find an h that agrees with f on training set
 - h is **consistent** if it agrees with f on all examples
- Finding a consistent hypothesis is not always possible
 - Insufficient hypothesis space:
 - E.g., it is not possible to learn exactly $f(x) = ax + b + x\sin(x)$ when $H =$ space of polynomials of finite degree
 - Noisy data
 - E.g., in weather prediction, identical conditions may lead to rainy and sunny days

لو انت مش قادر تلاقي general consistent hypothesis .. غالباً هتقول إن المشكله unrealizable .. والكلمه unrealizable .. هي مش وحشه يعني .. في ال practice .. هي بتقولك بس تعال نلاقي hypothesis space كبير جداً .. عشان نقدر نلاقي as many functions as possible .. واحنا في الآخر فعلاً بنزود احتمالية اننا نلاقي ال function اللي عاوزنها .. في trade off ما بين expressiveness و complexity .. لو عندك large hypothesis space .. هو expressive أكثر ويمكن فعلاً تلاقي hypothesis كويسه جداً .. بس انت محتاج تدور .. و التدوير بياخد وقت و data كتيره .. وبالتالي فلوس كتيره .. فتهحتاج داتا كتيره عشان انت هتعود تقارن hypothesis كتيره .. فمش دايماً شرط اني أخش علي طول علي largest hypothesis space .. بسبب ال tradeoff ده .. احنا بنحتاج بيقا عندنا space of reasonable functions .. وتشتغل

Inductive Learning

- A learning problem is **realizable** if the hypothesis space contains the true function otherwise it is **unrealizable**.
 - Difficult to determine whether a learning problem is realizable since the true function is not known
- It is possible to use a very large hypothesis space
 - For example: H = class of all Turing machines
- But there is a **tradeoff** between **expressiveness** of a hypothesis class and the **complexity** of finding a good hypothesis

دلوقت احنا جاهزين نتكلم علي أول machine learning technique

أول حاجه هنعرّف ال nearest neighbor classifier ب الفانكشن اللي اسمها $h(x)$.. ديه بتاخد انبوت $query(x)$.. و هتقوم مرجعه label y .. وده هيبقا ال label بتاع أقرب نقطه x^* .. فهنا x^* هي نقطه ليها أصغر مسافه .. والمسافه بنقيسها بطريقه معينه زي اللي في السلايد كدا ..

Nearest Neighbour Classification

- Classification function

$$h(x) = y_{x^*}$$

where y_{x^*} is the label associated with the nearest neighbour

$$x^* = \operatorname{argmin}_{x'} d(x, x')$$

- Distance measures: $d(x, x')$

$$L_1: d(x, x') = \sum_j^M |x_j - x'_j|$$

$$L_2: d(x, x') = \left(\sum_j^M |x_j - x'_j|^2 \right)^{1/2}$$

...

$$L_p: d(x, x') = \left(\sum_j^M |x_j - x'_j|^p \right)^{1/p}$$

$$\text{Weighted dimensions: } d(x, x') = \left(\sum_j^M c_j |x_j - x'_j|^p \right)^{1/p}$$

في حاجه اسمها L_1 و L_2 .. L_2 هي ال Euclidian distance .. سؤال حلو ... ازاى نفهم ايه اللي بيحصل في حالة ال classification .. انت بتقسم ال space .. انت عندك داتا بوينتس .. كل نقطه هي داتا بوينت .. و انا بشغل ال nearest neighbor classifier . انت تقدر تتخيل ان في شوية regions عندك .. جواهم شوية نقط و أي نقطه انت عاوز ت query هت map to the closest point . وبعدين ترجع ال Label بتاع النقطه ديه .. هو بيقسم ال input space لشوية مناطق . وال label بيبقا لكل منطقه .. ال partitioning معروف باسم Voronoi diagram .. وده بنجيبه عن طريق انك بتعمل nearest neighbor باستخدام ال nearest neighbor .. المناطق هتظهر عن طريق ان عندك boundaries علي هيئة straight lines .. هتبقا في النص علي طول ما بين ال 2 ريجنز ...

Voronoi Diagram

- Partition implied by nearest neighbor fn h
 - Assuming Euclidean distance



University of Waterloo

CS480/680 Spring 2019 Pascal Poupart

12

لما بتيجي تشتغل مع ال nearest neighbor هو unstable .. ليه .. لأن الفكره مبنيه علي انك ترجع ال class بتاع النقطه اللي اقرب حاجه لل query بتاعك .. بس الفكره هنا ايه اللي يحصل لو انت عندك شوية noise .. في الداتا اللي عندك .. الداتا بتيجي من شوية سينسورز ... والسينسورز بيبقا ليها ايرور في القرايات .. فبالنتالي اللي انت بتفكره اقرب حاجه ليك .. هو مش اقرب حاجه ليك .. او ممكن يكون فعلاً اقرب حاجه ليك .. بس هو متسمي غلط .. فالسؤال الجامد جداً هو ازاى تخلي ال nearest neighbor بيبقا more robust to noise .. أول فكره هو انك ت consider K-Nearest neighbor .. بدل ما انت بت produce a label based on a single nearest neighbor .. خلينا ن consider شوية neighbors و نبص علي most frequent class for that set of neighbors وبعدين هب نرجعهم .. فيكدا حتي لو نقطه واحده مكنش ليها أكبروسي كويس او كانت بعيدة شويه .. انت دلوقت معاك ال set اللي هترجعك الأقرب للصح فعلياً ..

K-Nearest Neighbour

- Nearest neighbour often instable (noise)
- Idea: assign most frequent label among k -nearest neighbours
 - Let $knn(x)$ be the k -nearest neighbours of x according to distance d
 - Label: $y_x \leftarrow mode(\{y_{x'} | x' \in knn(x)\})$

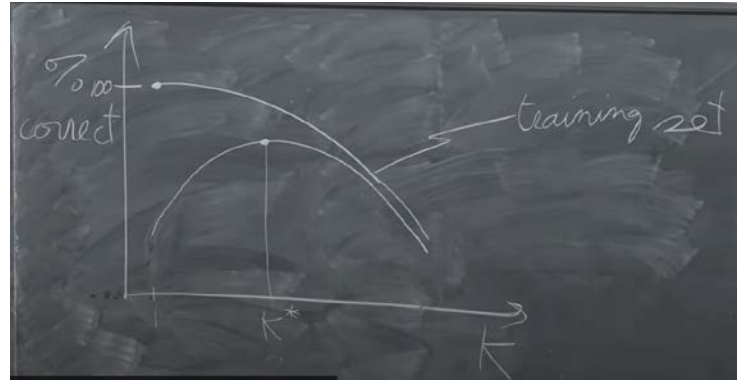
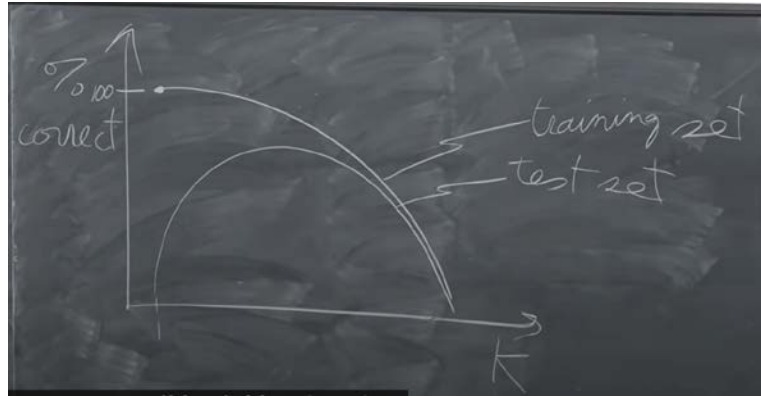
University of Waterloo

CS480/680 Spring 2019 Pascal Poupart

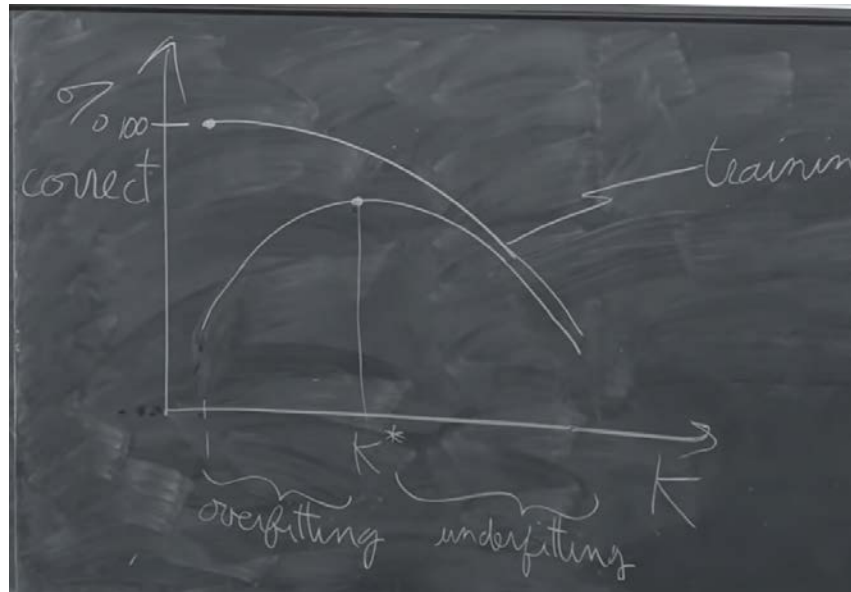
13

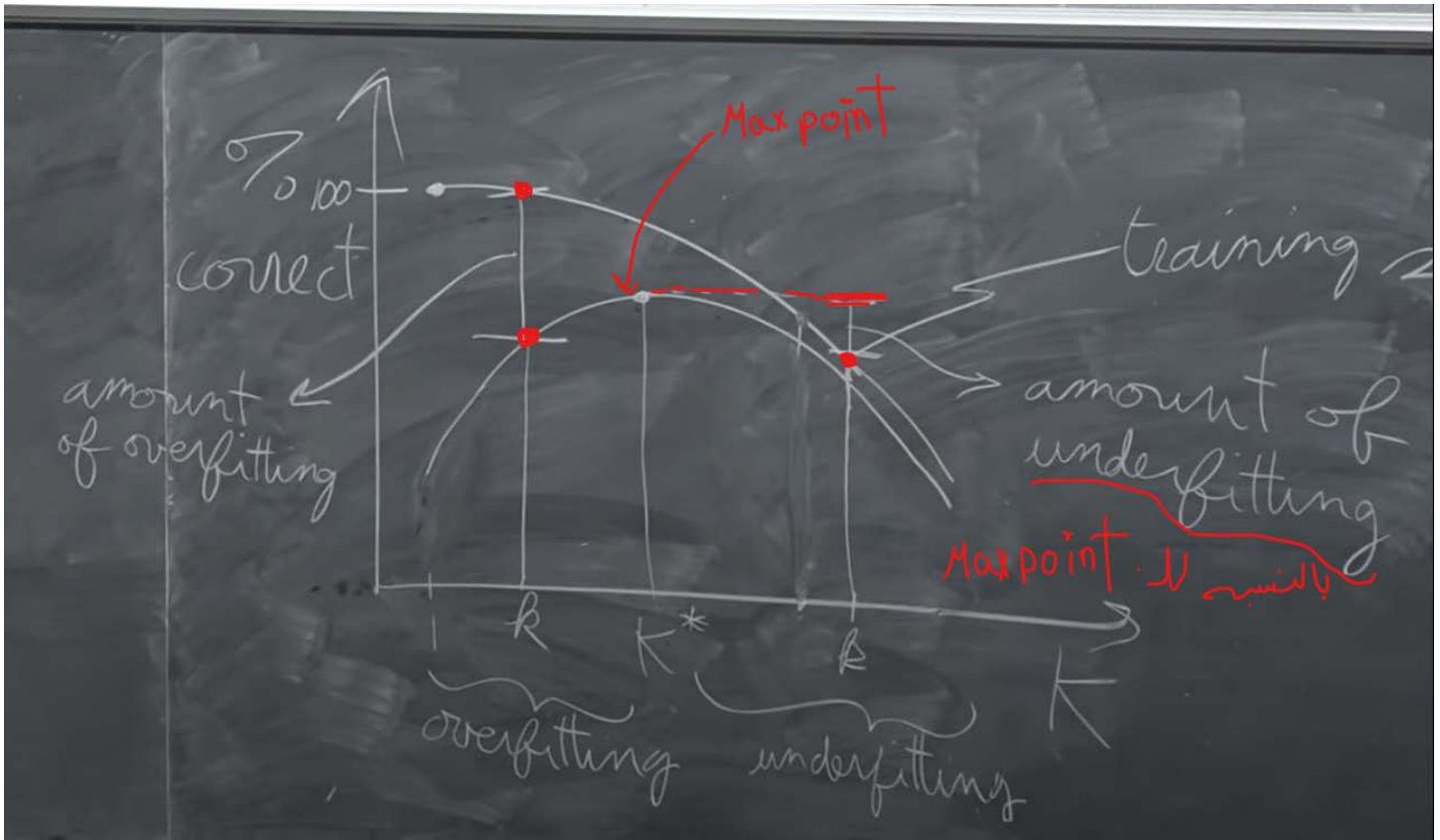
دلوقت لو انت هت consider K Nearest neighbor .. ايه هو الرقم الصحيح لل neighbors اللي عندك؟ .. ده مش حاجه سهله انك تقولها يعني بس جرب هي بتيجي بالتجربه ... في tradeoff ... لو انت كترت ال neighbors او انت قللتهم أوي .. هتلبس . اختار في النص .. تعال نبص علي الكيرف اللي الدكتور هيرسمه .. لو انت قلت ان عندك $K = 1$... هتلاقي عندك 100% كلهم صح .. ليه؟ .. عشان كل query هيبقا نقطه لما تلاقي ال nearest neighbor هتلاقي النقطه ديه ثاني .. وتدور علي ال label هتلاقيه correct label .. وترجعه فدايماً هترجع ال correct label ..

لو انت قلت عند $K = 1$.. وانت بت evaluate with respect to the training set .. انت عندك تريننج سيت .. وبعد ما عملت تريننج علي التريننج سيت .. جيت تبص علي ال accuracy بتاعت ال training هتلاقيها 100% .. ليه .. عشان كل query عباره عن نقطه في التريننج سيت .. اللي هو جزء من الداتا بيز بتاعتنا .. فلما هتيجي تلاقي ال nearest neighbor هتلاقي النقطه ديه ثاني .. فهترجع ال label بتاعها فهتلاقيك جيت ال 100% لكل الداتا ... only 1 neighbor الفكره في كيرف كذا لل K .. كل ما ال K زادت .. كل ما التريننج سيت accuracy قلت ... و العكس بالنسبه لل test ...set



الدكتور اتكلم بشكل بسيط علي ال underfitting .. اللي هو لو في hypothesis طلعتك accuracy سيئه جداً .. فأكد في hypothesis ثانيه بتطلع accuracy كويسه .. ال overfitting .. انت بت fit الداتا جامد جداً .. اللي هو كمان انك بت fit النويز اللي في الداتا جامد جداً برضو .. فانت بتضرب ال test accuracy .. في السلايد هنا بيقلوك ال futureAccuracy هنا هي الداتا الي ممكن تجيلك قدام اللي هي ال test accuracy .. الدكتور رجع للكيف اللي فوق ده ثاني عشان يوضح عليه ال overfitting و ال underfitting ...





Choosing K

- How should we choose K?
 - Ideally: select K with highest future accuracy
 - Alternative: select K with highest test accuracy
- **Problem:** since we are choosing K based on the test set, the test set effectively becomes part of the training set when optimizing K. Hence, we cannot trust anymore the test set accuracy to be representative of future accuracy.
- **Solution:** split data into **training, validation and test sets**
 - **Training set:** compute nearest neighbour
 - **Validation set:** optimize hyperparameters such as K
 - **Test set:** measure performance

سؤال جامد ثاني .. ازاي نختار K ... لو انت اخترت ال k بناءً علي ال test accuracy .. انت كدا بتعمل تريننج علي التست سيت مش علي التريننج سيت .. لأنك كدا كإنك effectively بتعمل training on test set .. فال K هنا هو بارمتر إننا بتعملو optimization .. وبكدا انت بتعملو optimization باستخدام ال test set .. وبكدا انت كإنك واخذ في الاعتبار ال test set كإنها في التريننج داتا .. او مال نعمل ايه .. split into 3 sets .. ترين و فاليديت و تست .. ال validation هتساعدك انك ت optimize the hyper parameter .. وهو ده الهدف من ال validation set .. تعال نبص علي الصودو كود مش الصودو كريم بتاع ال انك تختار ال K

Choosing K based on Validation Set

Let k be the number of neighbours

For $k = 1$ to max # of neighbours

$h_k \leftarrow \text{train}(k, \text{trainingData})$

$\text{accuracy}_k \leftarrow \text{test}(h_k, \text{validationData})$

$k^* \leftarrow \text{argmax}_k \text{accuracy}_k$

$h \leftarrow \text{train}(k^*, \text{trainingData} \cup \text{validationData})$

$\text{accuracy} \leftarrow \text{test}(h, \text{testData})$

Return $k^*, h, \text{accuracy}$

لو انا اخترت ال K بناءً علي how well our hypothesis performs with respect to to our validation set نويز .. ممكن هوب دابل كيك حظ كذا .. تلاقي واحده هي احسن قيمه لل K .. بس هو حظ عشان كان في شوية نويز علي الداتا بتاعت ال validation set ازاي نخلي الاختيار more robust ؟ .. انت تقدر تزود ال size of validation set .. ده هيجل المشكله بس هيطهرلك مشكله ثانيه .. انك هتقلل الداتا بتاعت التريننج .. فاتضح إن في طريقه تقدر تستخدم ال training and validation set اللاتنين للاتنين ال training and validation و ده اللي هو ال cross validation .. ده بيحصل عن طريق انك تاخذ الداتا . وتقسمها نصين .. وتفضل تقسم تقسم تقسم تقسم ... وبعدين تاخذ ال average K-fold cross .. accuracy across all the experiments with respect all those splits .. validation

Robust validation

- How can we ensure that validation accuracy is representative of future accuracy?
 - Validation accuracy becomes more reliable as we increase the size of the validation set
 - However, this reduces the amount of data left for training
- Popular solution: **cross-validation**

Cross-Validation

- Repeatedly split training data in two parts, one for training and one for validation. Report the average validation accuracy.
- **k-fold cross validation**: split training data in k equal size subsets. Run k experiments, each time validating on one subset and training on the remaining subsets. Compute the average validation accuracy of the k experiments.
- Picture:

تعال نرسم صورته تكون أوضح شويه للمعني ...

$k=4$ 4-fold cross validation

