

Q&A

Q: Why do we shuffle the examples in SGD?

A: This is how we do sampling *without* replacement

1. **Theoretically** we can show sampling **without replacement** is not significantly worse than sampling with replacement (Shamir, 2016)
2. **Practically** sampling without replacement tends to work better

Q: What is “bias”?

A: That depends. The word “bias” shows up all over machine learning! Watch out...

1. The additive term in a linear model (i.e. b in $w^T x + b$)
2. Inductive bias is the principle by which a learning algorithm generalizes to unseen examples
3. Bias of a model in a societal sense may refer to racial, socio-economic, gender biases that exist in the predictions of your model
4. The difference between the expected predictions of your model and the ground truth (as in “bias-variance tradeoff”)

2

طيب بسم الله .. احنا عاوزين نفكر في السامبل كومبلكستي .. في آخر مره احنا وصلنا لأول نظريه

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	<p>Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.</p>	
Infinite $ \mathcal{H} $	<p>True error</p>	<p>Training error</p>

6

هنا انت عندك upper bound علي ال true error اللي انت متعرفهوش بناء علي ال training error الي احنا بنبقا عارفينو ... طيب الدكتور حاطط مثال وعاوزنا نفكر في النظرية كإنها مثال .. والمثال هنا هو هيبقا hypotheses class of conjunctions .. فاحنا هنقول ان H is a class of conjunctions over vectors X .. والفكتورز إكس دول هيبقا اصصفار ووحدات ..

Example: Conjunctions

Question:

Suppose H = class of conjunctions over x in $\{0,1\}^M$

Example hypotheses:

$$h(x) = x_1 (1-x_3) x_5$$

$$h(x) = x_1 (1-x_2) x_4 (1-x_5)$$

If $M = 10$, $\epsilon = 0.1$, $\delta = 0.01$, how many examples suffice according to Theorem 1?

$$|H| = 3^M$$

Answer:

- A. $10 * (2 * \ln(10) + \ln(100)) \approx 92$
- B. $10 * (3 * \ln(10) + \ln(100)) \approx 116$
- C. $10 * (10 * \ln(2) + \ln(100)) \approx 116$
- D. $10 * (10 * \ln(3) + \ln(100)) \approx 156$
- E. $100 * (2 * \ln(10) + \ln(10)) \approx 691$
- F. $100 * (3 * \ln(10) + \ln(10)) \approx 922$
- G. $100 * (10 * \ln(2) + \ln(10)) \approx 924$
- H. $100 * (10 * \ln(3) + \ln(10)) \approx 1329$

$$\frac{1}{0.1} \left[\ln(3) + \ln(100) \right]$$

طبيب الدكتور بيقول في المثال ده اننا عاوزين نجيب عدد الامثله اللي هتبقا كفايه عشان تحقق theorem 1 .. الدكتور اتكلم علي ان المشكله كلها في حساب ال size of H .. اللي هو في النظرية عندك $|H|$ المهم هنا الدكتور قال هو الحجم اصلا عبارته عن ايه ... عبارته عن all the possible conjunction function you can come up with ... الدكتور كمل كلامه وقال ان في 3 اوبشنز لكل فاريل من ال M variables اللي عندك .. كل فاريل هنا هيبقا ليه 3 حالات . يا إما included او included but negated او excluded all together .. من $x_1, x_2, x_3, \dots, x_{10}$.. كل واحد ليه 3 اختيارات زي دول .. فبالتالي الحجم بتاع ال hypothesis space هيبقا 3^{10} ... فبالتالي الاجابه الصح هتبقا D ... طبيب الدكتور بيقول ده معناه ايه المثال ده .. معناه انك لو عندك 156 examples فبالتالي according to theorem 1 .. ده هيبقا كفايه عشان باحتمالية 99% كل ال hypotheses اللي هو all conjunction functions هيبقا ليهم صفر تريننج ايرور علي الداتا سيت بتاعتك اللي هي 156 مثال ... وهيبقا ليهم true error اقل من 10% الحاجه الحلوه هنا ان علي ارض الواقع احنا مش بنستخدم النظرية ديه عشان نحدد عدد الامثله اللي هحتاجها عشان نجيب البرفورمنس اللي احنا عاوزينو .. هي بس مجرد بتديك sense كذا انك تقول انت محتاج كام مثال يعني انت ممكن تحتاجهم عشان تجيب برفورمنس كويس تحت شوية assumptions قلناهم المره اللي فانت تحت ال PAC model ...

طبيب الحاجه الثانيه اللي هتقول عليها هي ال Realizable vs Agnostic case .. ال Relalizable setting هنا انت عندك ال c star جوا ال hypothesis space ... انما ال Agnostic setting ... هي ان برضو ال c star in the hypotheses space أو انها مش في ال hypotheses space احنا مش متأكدين ... "عشان نبقا ناس كوريكتور" ..

طبيب السؤال هنا عشان نخلي نظريه رقم 1 ديه ليها قيمه يعني احنا فعلا بنحتاج نجيب قيم الإيسيلون و قيم الدلتا ... هنا انت بتختارهم بايدك ... طبيب تعال نرجع للفرق ما بين ال Realizable vs Agnostic ... طبيب اللي عاوزين نقولو هنا ان ال Agnostic case بتيجي لما انت بتفكر في ان أيا كان الفانكشن اللي عملت جينيريشن للداتا .. هي مش معاك في الهيبوثيسيز سبيس ... فال H finite هنا احنا لو فكرنا مثلا زي في ال class of decision trees ... مثلا لو عندك ماكس ديبث 10 .. ليه H clear finite تقدر تبص علي كل الهيبوثيزيز اللي عندك ... بس فعليا كل الداتا سيتس اللي هتعمل تريننج عليها .. هتبقا في ال Agnoistic ليه .. عشان اللي عمل ال labeling هو انسان .. مش شجره .. "انا شجره" .. فاللي بنقولو في ال Agnostic .. في السلايد الجايه

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable $\epsilon \in \mathcal{H}$	Agnostic $\epsilon \notin \mathcal{H}$
Finite $ \mathcal{H} $	Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.
Infinite $ \mathcal{H} $		True Train

طيب اللي عاوزين نلاحظ هنا هو الاختلاف بتاع ال utility بتاع ال 2 باوندر اللي عندنا ...

	Realizable	Agnostic
Finite $ \mathcal{H} $	<p>1. Bound is inversely linear in epsilon (e.g. halving the error requires <u>double the examples</u>)</p> <p>2. Bound is only logarithmic in \mathcal{H} (e.g. quadrupling the hypothesis space only requires double the examples)</p>	<p>1. Bound is inversely quadratic in epsilon (e.g. halving the error requires <u>4x the examples</u>)</p> <p>2. Bound is only logarithmic in \mathcal{H} (i.e. same as Realizable case)</p>
Infinite $ \mathcal{H} $	<p>Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.</p> <p><i>i.e. $\mathcal{H} = 3^M$ $\log = M \log 3$</i></p>	<p>Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $R(h) - \hat{R}(h) \leq \epsilon$.</p> <p><i>$\epsilon \rightarrow 0.01 \rightarrow N$ $\epsilon' \rightarrow 0.005 \rightarrow 4N$ أربع المبرك</i></p>

9

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	<p>Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.</p>	<p>Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $R(h) - \hat{R}(h) \leq \epsilon$.</p>
Infinite $ \mathcal{H} $		

We need a new definition of "complexity" for a Hypothesis space for these results (see VC Dimension)

12

طيب لو عندنا linear decision boundary عندنا انفتت نامير اوف هايپوثيزيز ..

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

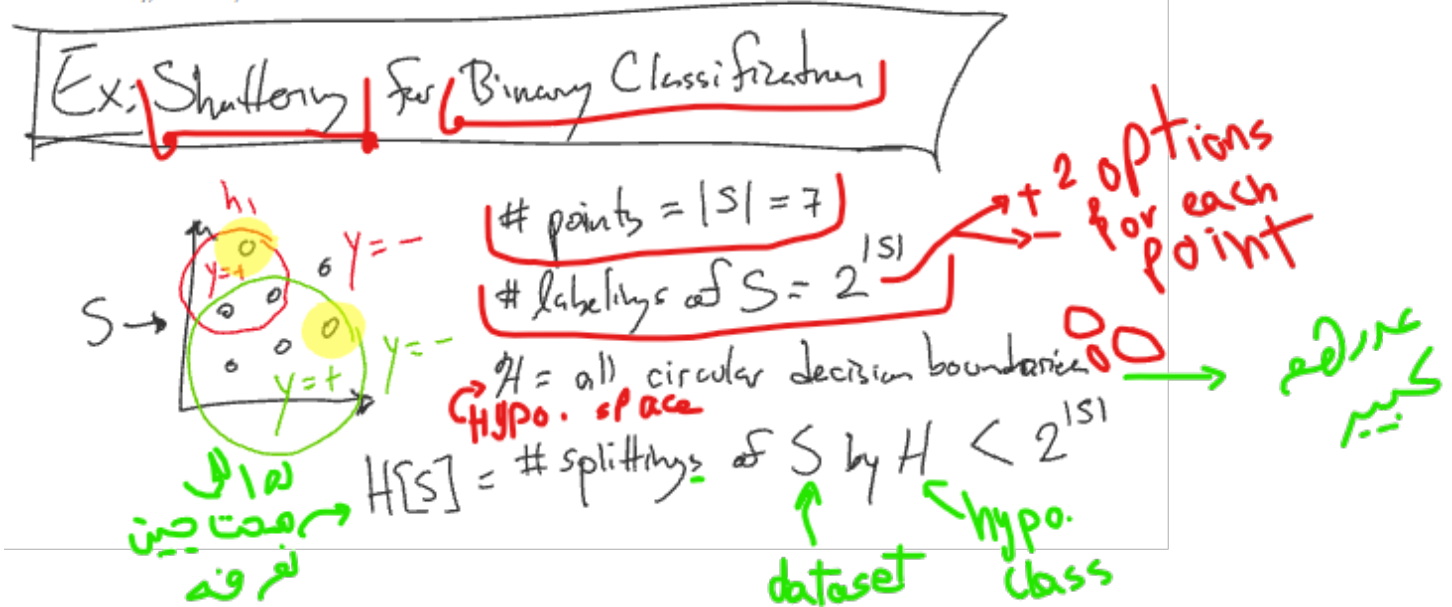
Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	<p>Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.</p>	<p>Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $R(h) - \hat{R}(h) \leq \epsilon$.</p>
Infinite $ \mathcal{H} $	<p>Thm. 3 $N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.</p>	<p>Thm. 4 $N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $R(h) - \hat{R}(h) \leq \epsilon$.</p>

13

PAC Learning

Wednesday, March 18, 2020 11:54 AM



الدكتور بيقول ان في labelings .. السيركيولار ديسيजन باوندري مبيقدرش ي replicate .. يعني مثلا لو خدنا النقط اللي متعلم عليها بالاصفر فوق .. لو الاتنين كانو +ve label .. انت معندكش circular decision boundary اللي بيقدر ي label the dataset بالطريقه ديه .. في الحاله ديه عدد ال splitting هيبقا اكيد اكيد اقل من $2^{|S|}$.. وده اللي احنا وصلناو في السلايدز ..

لو جينا نفكر في التعريف ده بيقول ايه .. هو بيقول ان ال hypotheses space بي shatter S لو عدد ال splitting المتاح ليك هو ببساوي $2^{|S|}$... فال circular decision boundary للنقط الصفرا مقدرش انو ي shatter that particular dataset S ..

طبيب هنعرف ايه هو ال VC-dimension .. ال VC-dimension بتاع ال hypotheses space هو حجم أكبر dataset S تقدر تقسمها باستخدام ال H .. السؤال هنا ازاي نجيب VC-dimension بتاع some particular hypotheses space ..

Shattering, VC-dimension

Definition:

$H[S]$ - the set of splittings of dataset S using concepts from H .

H shatters S if $|H[S]| = 2^{|S|}$.

A set of points S is shattered by H if there are hypotheses in H that split S in all of the $2^{|S|}$ possible ways; i.e., all possible ways of classifying points in S are achievable using concepts in H .

Definition: VC-dimension (Vapnik-Chervonenkis dimension)

The **VC-dimension** of a hypothesis space H is the cardinality of the largest set S that can be shattered by H . size of

If arbitrarily large finite sets can be shattered by H , then $VCdim(H) = \infty$

Ex: VC Dimension for Linear Sep.

\mathcal{H} = linear separators in 2D $M=2$

To prove $VC(\mathcal{H}) = d = M+1 = 3$

There exists set S

① $\exists S \subseteq \mathcal{X}$ s.t. $|S|=d$ and \mathcal{H} shatters S

② $\nexists S \subseteq \mathcal{X}$ s.t. $|S|=d+1$ and \mathcal{H} shatters S

not exist
sub set of all possible points

قدرة اننا
نقسم النقط
في كل اماكن
possible ways

فعلشان تثبت ال VC-Dimension اللي انت محتاج تقول انو موجود هو example dataset ال \mathcal{H} تقدر ت shatter .. وديه أكبر داتا سيت ال \mathcal{H} can shatter .. "انا مش فاهم" ... أول سؤال ازاى نثبت رقم 1 في الصورة اللي فوق:

To show ①:

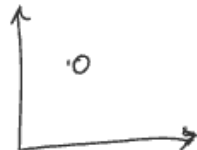
size of S

Pick one dataset (unlabeled) for d

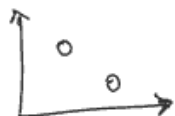
List

List all possible labelings of S $\rightarrow +ve$ $\rightarrow -ve$ Show that you can shatter S

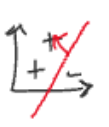
$d=1$



$d=2$



$d=3$

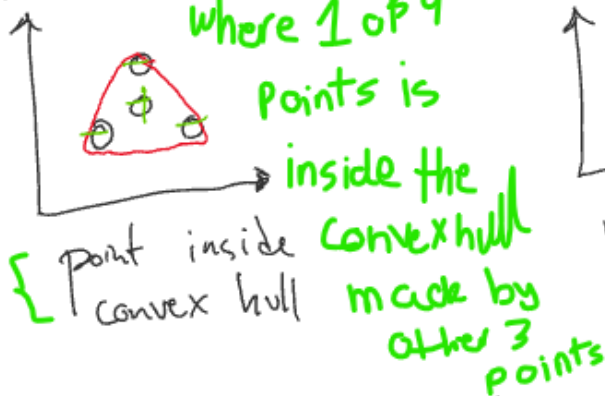


اللي عاوزين نعملو ناو هو نثبت رقم 2:

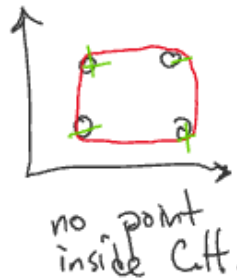
To show (2):

Divide all datasets into two categories

Case A all datasets where 1 of 4



Case B



For all S of these two cases, H cannot shatter S

اللي بنقولو هنا هو ان لو عندنا مثلا في حالة A .. لو انت تخيلت ان عندك داتا سبت من الفورم اللي بالشكل ده .. و البوينتس كانت ان اللي في النص بس هو الموجب و كل اللي الكورنرز كانوا سالب ... في الحاله ديه مفيش linear decision boundary بيقدري correctly classify that .. ومش مهم ال plus point ديه فين جوا ال convex hull او ايه هو الكونفجریشن بتاع ال 3 نقط التانيين .. فأي دات سبت هتجيبها لي مش هيبقا في linear decision .. boundary

الدكتور هيحط سايد نوت عبارته عن سؤال ..

Side Note

Q: What is dataset of size $d=3$ that can't be shattered?



طيب اللي علوزين نوضحو ان في فرق ما بين ال distinction between there exists and for all .. في ال VC-dimension .. لما بنيجي نثبت ال VC-Deimension بيحتاجنا نثبت ان في داتا سبت حجمها d تقدر ت شاتر و كمان انك تثبت ان مفيش داتا سبت حجمها $d+1$ منقدرش شاتر .. الدكتور بس حب يوضح السلايد الجايه ... ان ال VC بتعتمد علي exists انما ال shattering بيعتمد علي ال for all ..

∃ vs. ∀

VCDim

- Proving **VC Dimension** requires us to show that **there exists** (\exists) a dataset of size d that can be shattered and that **there does not exist** (\nexists) a dataset of size $d+1$ that can be shattered

Shattering

- Proving that a particular dataset can be **shattered** requires us to show that **for all** (\forall) labelings of the dataset, our hypothesis class contains a hypothesis that can correctly classify it

طيب في امثله ثانيه:

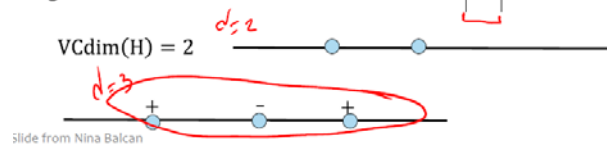
Shattering, VC-dimension

If the VC-dimension is d , that means **there exists** a set of d points that can be shattered, but there is **no** set of $d+1$ points that can be shattered.

E.g., $H =$ Thresholds on the real line



E.g., $H =$ Intervals on the real line



slide from Nina Balcan

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	Thm. 1 $N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 2 $N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.
Infinite $ \mathcal{H} $	Thm. 3 $N = O(\frac{1}{\epsilon} [VC(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	Thm. 4 $N = O(\frac{1}{\epsilon^2} [VC(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.

29

طبيب الدكتور بيقول احنا لحد دلوقت بنتكلم علي عدد الأمثله اللي هحتاجها .. دلوقت هنخش في اتجاه مختلف شويه هو ال slt اللي هو statistical learning theory ... فالدكتور في السلايد الجايه ديه كان بيشرح ازاي نروح من نظريه رقم 1 ل 1 corollary ... انك هتحل للإيسيلون .. ده بيحط upper bound علي التروو إيرور ل fixed number of examples

SLT-style Corollaries

Thm. 1 $\frac{1}{N} \geq \frac{1}{N} \left[\log(|\mathcal{H}|) + \log\left(\frac{1}{\delta}\right) \right]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.

Solve the inequality in Thm.1 for epsilon to obtain Corollary 1

Corollary 1 (Realizable, Finite $|\mathcal{H}|$). For some $\delta > 0$, with probability at least $(1 - \delta)$, for any h in \mathcal{H} consistent with the training data (i.e. $\hat{R}(h) = 0$),

$$R(h) \leq \frac{1}{N} \left[\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right]$$

We can obtain similar corollaries for each of the theorems...

SLT-style Corollaries

Corollary 1 (Realizable, Finite $|\mathcal{H}|$). For some $\delta > 0$, with probability at least $(1 - \delta)$, for any h in \mathcal{H} consistent with the training data (i.e. $\hat{R}(h) = 0$),

$$R(h) \leq \frac{1}{N} \left[\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right]$$

Corollary 2 (Agnostic, Finite $|\mathcal{H}|$). For some $\delta > 0$, with probability at least $(1 - \delta)$, for all hypotheses h in \mathcal{H} ,

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2N} \left[\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right]}$$

SLT-style Corollaries

Corollary 3 (Realizable, Infinite $|\mathcal{H}|$). For some $\delta > 0$, with probability at least $(1 - \delta)$, for any hypothesis h in \mathcal{H} consistent with the data (i.e. with $\hat{R}(h) = 0$),

$$R(h) \leq O \left(\frac{1}{N} \left[\text{VC}(\mathcal{H}) \ln \left(\frac{N}{\text{VC}(\mathcal{H})} \right) + \ln \left(\frac{1}{\delta} \right) \right] \right) \quad (1)$$

Corollary 4 (Agnostic, Infinite $|\mathcal{H}|$). For some $\delta > 0$, with probability at least $(1 - \delta)$, for all hypotheses h in \mathcal{H} ,

$$R(h) \leq \hat{R}(h) + O \left(\sqrt{\frac{1}{N} \left[\text{VC}(\mathcal{H}) + \ln \left(\frac{1}{\delta} \right) \right]} \right) \quad (2)$$

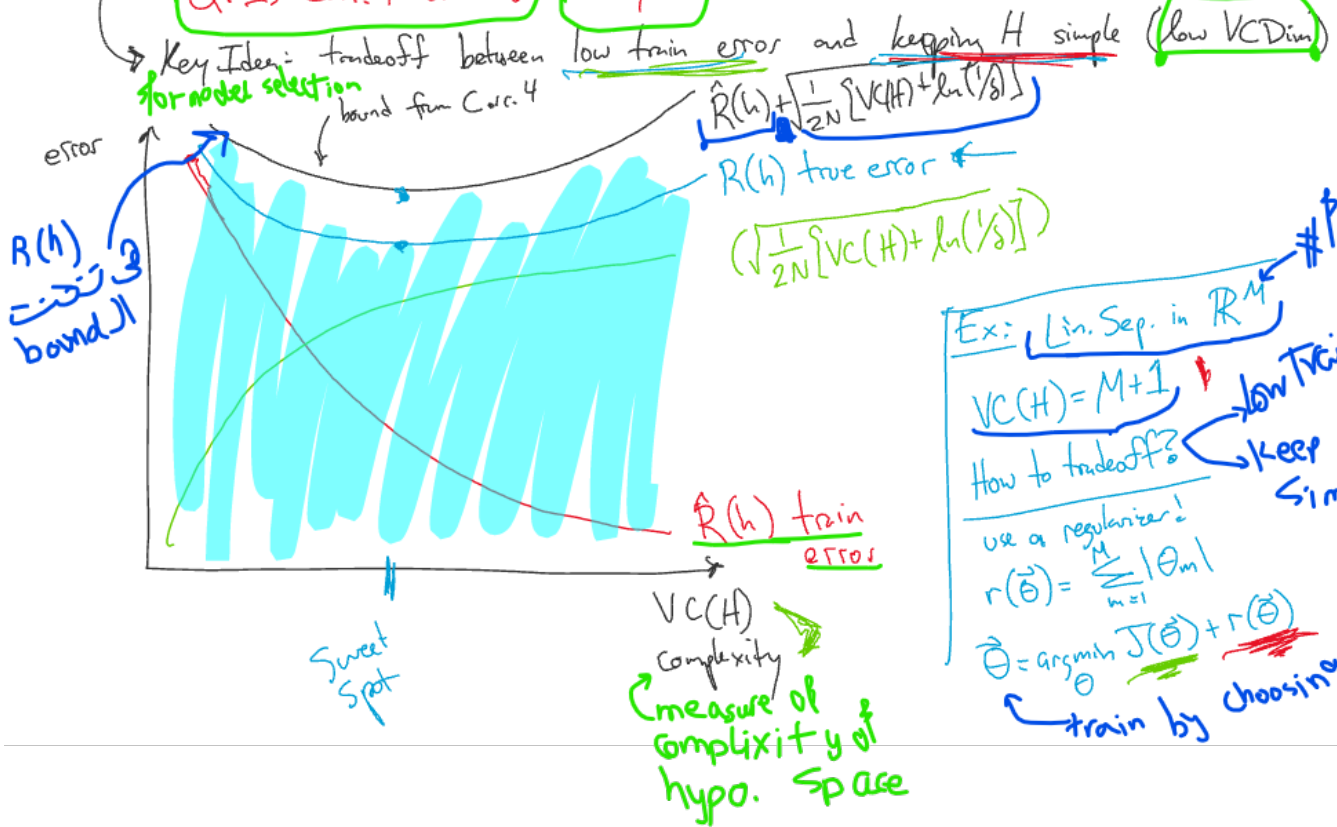
طبيب السؤال المهم .. ازاي هنعمل model selection :

.. empirical risk minimization هو انت عاوز تفكر في

Model Selection

Q. Is Corc. 4 useful?

A: Yes!



Recipe for Machine Learning

- ① Given data $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$
- ② a) Choose a decision function $h_{\vec{\theta}}(\vec{x}) = \dots$ parameterized by $\vec{\theta}$
- b) Choose an objective function $J_D(\vec{\theta}) = \dots$ relies on data
- ③ Learn by choosing parameters that optimize the objective $J(\vec{\theta})$
 $\hat{\vec{\theta}} \approx \arg\min_{\vec{\theta}} J_D(\vec{\theta})$
- ④ Predict on new test sample \vec{x}_{new} using h
 $\hat{y} = h_{\vec{\theta}}(\vec{x}_{new})$

Decision Functions

- Perceptron: $h_{\theta}(\vec{x}) = \text{sign}(\vec{\theta}^T \vec{x})$

- Linear Regression: $h_{\theta}(\vec{x}) = \vec{\theta}^T \vec{x}$

- Discriminative Models: $h_{\theta}(\vec{x}) = \underset{y}{\text{argmax}} p(y|\vec{x})$

Log. Reg. $p(y|\vec{x}) = \sigma(\vec{\theta}^T \vec{x})$

- Generative Models: $h_{\theta}(\vec{x}) = \underset{y}{\text{argmax}} p_{\theta}(\vec{x}, y) \rightarrow \text{joint distr.}$

Naive Bayes $p(\vec{x}, y) = p(y) \prod_{n=1}^M p(x_n | y)$

- NN for classification: $p(y|\vec{x}) = \sigma(W^{(2)T} \sigma(W^{(1)T} \vec{x} + b^{(1)}) + b^{(2)})$

Obj. Function

① MLE: $J(\theta) = - \sum_{i=1}^N \log p(x^{(i)}, y^{(i)})$

Conditional → ② MLE: $J(\theta) = - \sum_{i=1}^N \log p(y^{(i)} | x^{(i)})$

③ L2 Reg: $J'(\theta) = J(\theta) + \lambda \|\theta\|_2^2$

④ L1 Reg: " " + $\lambda \|\theta\|_1$

Opt Method

- Gradient Descent

$$\theta \leftarrow \theta - \gamma \nabla J(\theta)$$

- SGD where

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta)$$

$$\theta \leftarrow \theta - \gamma \nabla J^{(i)}(\theta) \text{ for } i \sim \text{Uniform}(\{1, \dots, N\})$$

- mini-batch SGD

- Closed Form

→ Partial

① compute partial derivatives

② set to zero and solve