

احنا كنا بنحاول ن classify 3 dimensional points of different animals .. الحيوانات كانوا polar bear, sea lion and whales .. و في الآخر قلنا اننا عندنا موديل اللي هو ال multinomial logistic regression و هنا قلنا انو متعرف علي اساس انو probability distribution .. في حين ان ال probabilities جات من انك تبص علي ال signed distance of a point بالنسبة لواحد من ال 3 hyperplanes عندني .. و ال signed distance بمعنى انها لو علي الناحية البوزيتيف بيقا هي مسافه فعلاً .. انما لو كانت علي النيجاتيف سايد .. فانت بتضربها فسالب واحد ههه . عشان الديستانس دايماً بوزيتيف .. وبعدين قلنا نرفع الكلام ده ل exp عشان نجيب رقم موجب ما بين الصفر و المالا نهايه .. بعدها هت normalize الكلام ده .. قلنا بعد كذا ان زي ما كان بيحصل في ال binary logistic regression ... انت تقدر تفكر في الموديل علي اساس انو coin flipping distribution ... والواي جايه من Bernoulli distribution والبارمتر فاي بتاعه هو عبارته عن فانكشن في البارمترز ثيتا و الفيتشرز اكس .. هنا بنقول انك تقدر تعمل weighted die roll انت معندكش غير die واحده بس .. و الواحده ديه ليها k sides .. وكل ناحيه ليها الاحتمالية بتاعتها والاحتماليه ديه بتتبقا deterministic function of theta and X والواي في الحاله ديه هي بتتبقا شبه ال weighted die roll in a multinomial case .. احنا هنكمل علي المره اللي فانت .. دلوقت علوزين نفكر ال learning هنا معنا ايه وبيحصل ازاى .. فعشان نتعلم مع ال maximum likelihood estimation اللي بنحتاج نعملو اننا نكتب ال objective function و ديه بتتبقا عبارته عن conditional log-likelihood و هنا هنعرفها علي اساس انها فانكشن L(theta) و هي اللوج بتاعت ال log-likelihood و ديه بتتبقا ال multiplication from 1 to N ال p(y|x, theta) .. لما بندخل ال Log في النص بنأخذ summation مش multiplication .. ودلوقت هنعرف فانكشن اسمها J(theta) ... و هي ال -1/N L(theta) ... و الفانكشن اللي اسمها J(theta) ديه بتتبقا convex فانت تقدر تستخدم ال GD and SGD عشان ت minimize J(theta) .. فعشان تستخدم واحده من الاتنين اللي هم ال GD or SGD انت هتحتاج تحسب ال gradient فهتحتاج تجيب ال partial derivatives .. فهتقول ان اللي انت علوزو هنا هو ال partial derivative of Ji(theta) بص ع الصوره

③ Learning by MLE Conditional Log-Likelihood :

$$\ell(\theta) = \log \left[\prod_{i=1}^N p(y^{(i)} | x^{(i)}, \theta) \right] = \sum_{i=1}^N \log p(\dots)$$

$$J(\theta) = -\frac{1}{N} \ell(\theta) \quad \leftarrow \star \text{convex} \star$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) \quad \text{use GD or SGD}$$

Gradient :

$$\frac{\partial J^{(i)}(\theta)}{\partial \theta_{km}} = \frac{\partial}{\partial \theta_{km}} \left(-\log p(y^{(i)} | x^{(i)}, \theta) \right)$$

$$= \dots \quad \leftarrow \text{HW4} \quad \begin{cases} 1 & \text{if } y^{(i)} = k \\ 0 & \text{otherwise} \end{cases}$$

$$= - \left(\mathbb{I}(y^{(i)} = k) - p(y^{(i)} | x^{(i)}, \theta) \right) x_m^{(i)}$$

$$\frac{\partial J^{(i)}(\theta)}{\partial \theta_k} = \nabla_{\theta_k} J^{(i)}(\theta) = - \left(\mathbb{I}(y^{(i)} = k) - p(y^{(i)} | x^{(i)}, \theta) \right) \tilde{x}^{(i)}$$

$$= \begin{bmatrix} \frac{\partial J^{(i)}(\theta)}{\partial \theta_{k1}} \\ \frac{\partial J^{(i)}(\theta)}{\partial \theta_{k2}} \\ \vdots \\ \frac{\partial J^{(i)}(\theta)}{\partial \theta_{km}} \end{bmatrix} \quad \text{compute for each } k$$

ال || ديه اسمها binary indicator .. ديه كل الي بتعملو انها تقول 1 لو ال $y_i = k$.. غير كذا هتبقا بصفر .. المهم لم كل دول وحطهم في ال gradient قال gradient هنا هيبقا زي الصوره اللي فوق كذا بالظبط ... خد بالك بس من حاجه .. الجريدينت اللي بتحسبو ده اسمو vector derivative .. وهنتكلم عنو في ال NNS .. المهم انت مش لازم تعرف ازاى تتعامل معاهم يعني .. بس ديه مشتقه بالنسبة لفيكتور كامل ... فانت لما بتعمل كذا بيطلعلك الجريدينت .. وده الكولوم فيكتور اللي انت شايفو فوق في الصوره ديه ... النقطه اللي الدكتور حاططها ديه اللي هي ال J(theta) بس الدكتور مكسل يكتب كل شويه

يعني ... خذ بالك ان الفيكتور ده مش ال full gradient اللي انت مهتم بيها .. اللي احنا فعلاً مهتمين بيه هو ال full gradient for the entire matrix of parameters .. انت هتحتاج تحسب الفيكتور اللي انت شايفو بالازرق فوق ده لكل class k عندك .. وده هيديلك ال k rows of the gradient of the matrix .. أول ما هتعمل الكلام ده .. تعال نستخدم ال SGD او GD عشان نتعلم ال optimization وبعدها نعمل ال prediction لأحسن كلاس في الموديل بتاعنا .. احنا شغنا الكلام ده في حاجات كتير زي مثلاً majority vote او binary logistic regression و دلوقت ل multinomial logistic regression .. و دلوقت عاوزين نتكلم علي حته practical كذا .. في غلطه شائعته .. الكود اللي في الصورة الجايبه ديه فيه غلطه ..



Debug that Program!

In-Class Exercise: Think-Pair-Share

Debug the following program which is (incorrectly) attempting to run SGD for multinomial logistic regression

Buggy Program:

while not converged:

for i in shuffle([1,...,N]):

for k in [1,...,K]:

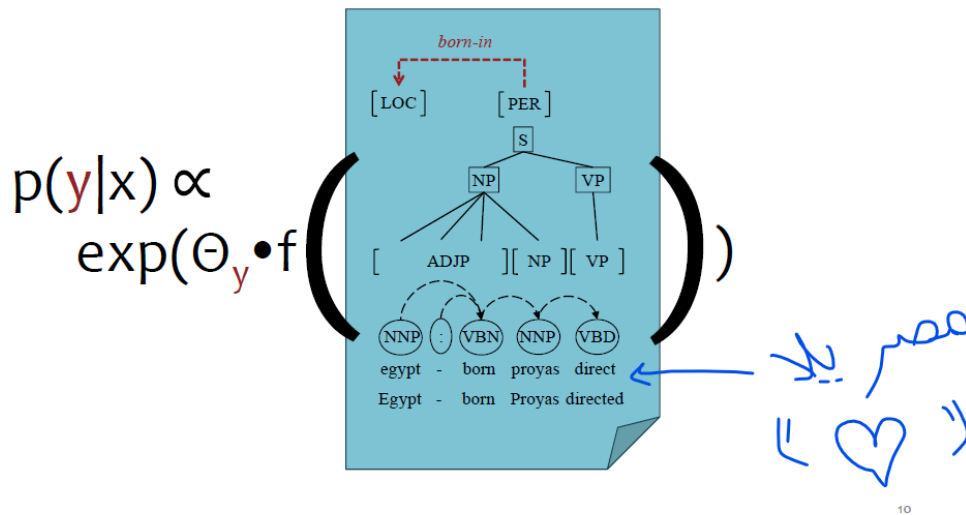
theta[k] = theta[k] - gamma * grad(x[i], y[i], theta, k)

Assume: grad(x[i], y[i], theta, k) returns the gradient of the negative log-likelihood of the training example (x[i], y[i]) with respect to vector theta[k]. gamma is the learning rate. N = # of examples. K = # of output classes. M = # of features. theta is a K by M matrix.

افتراض انك بتحاول ترن ال SGD وعندك 2 بارمترز theta 1 and theta 2 .. هل ال SGD or GD بتاخذ خطوه في ال direction بتاع ال steepest descent ولا الاول بتاخذ ال steepest descent لبارمتر و بعدين البارمتر الثاني .. الدكتور قال ان المفروض اننا نحسب الجريدينت كلو بتاع الثيتا وبعدين نأبدت الثيتا بالجريدينت ده .. انت فعلاً ممكن تفكر فيها كإنيك عندك مشكله فيها ال M = 1 and k = 2 .. كذا الثيتا ماتركس فيها 2 entries .. انت هنا المفروض تحسب ال full gradient matrix اللي هي دلوقت بقت فيكتور ليه 2 entries .. وبعدين تاخذ خطوه في الاتجاه ده .. اللي هو الخط الأحمر ... انما اللي انت بتعملو في الصورة الي فوق .. انك بتحسب ال gradient لواحد من الاتنين بارمترز دول .. اللي هو قيمه واحده من القيمتين بتوع ال K .. وبعدين تقوم واخذ خطوه في الاتجاه ده و تعمل update للثيتا .. وبعدين ت call the gradient ثاني .. للقيمه الثانيه بتاعت ال k وتقوم واخذ خطوه ثانيه ... وده ممكن يودينا لمكان ثاني .. فالامبلمنتيشن عشان يتصلح ... حط variable tmp في النص ما بين ال 2 loops .. الفور لوب ال اللي في النص ديه هتسبها كل قيم ال gradient .. وتقوم مستخدمها عشان تعمل update to theta ..

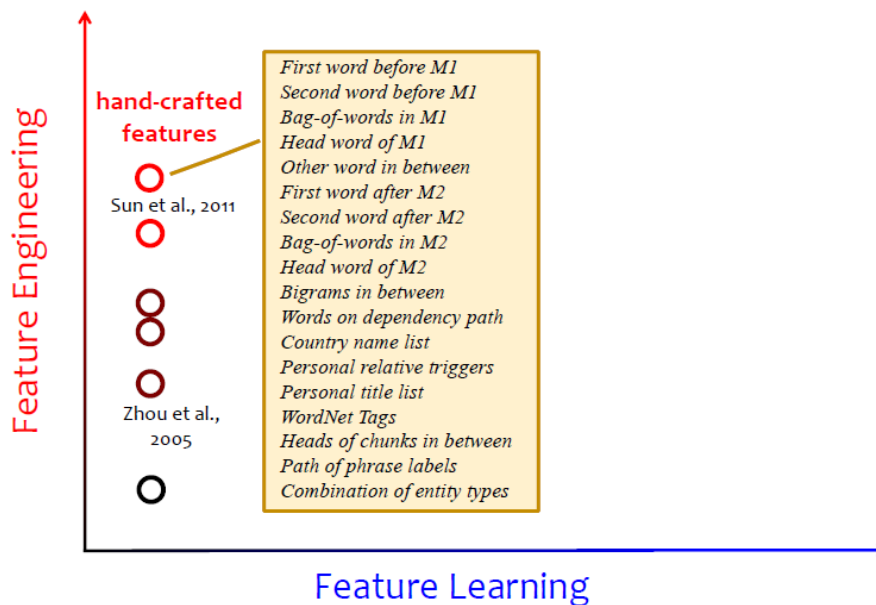
تعال دلوقت نفكر في مشكله مهمه .. احنا اتكلمنا علي الفيتشرز كتير ... بس الفكره بنجيبهم منين ... افترض ان عندك جمله ... جمله نفسها ممله .. انما ايه رأيك لو عدلنا علي الجملة مثلاً خليناها كلها lower-case او انك شلت شويه حاجات منها .. بعدها تقدر تضيف مثلاً جزء من speech tags للجمله .. او تضيفلهم حاجات كتيره .. بص عالصوره ... اول ما يكون عندك كل ده .. دخلو في الفانكشن f و ده هيرجعلك فيتشر فيكتور ريبزنتيشن لكل ال ستركتشر ده .. الفيكتور ده هيبقا عامل ازاي .. تعال نبص ع الجملة

Handcrafted Features



هو مين بروياس ده .. هل ده متجوز يعني .. هل ههو مولود في مصر .. اشتري الكلب بتاعه من مصر .. عمل ايه يعني .. هنا بقا الموديل هيقدر انو ي pick out the correct relationship ما بين الحاجات ديه باستخدام الفيبينشر فيكتور بتاعك .. ازاي الناس بتعمل كدا .. طريقه انهم بيعملو ريسيرش .. ازاي ن handcraft احسن فانكشن f ... بص ع الصورة .. انت هتكتب كود بيرن f deterministic function وببشده اهم الحاجات و يجيلك فيكتور طويل عريض كدا فيه كل الحاجات المهمه اللي انت جبتها .. زي مثلاً الحاجات اللي بالاصفر في الصورة اللي تحتينا ..

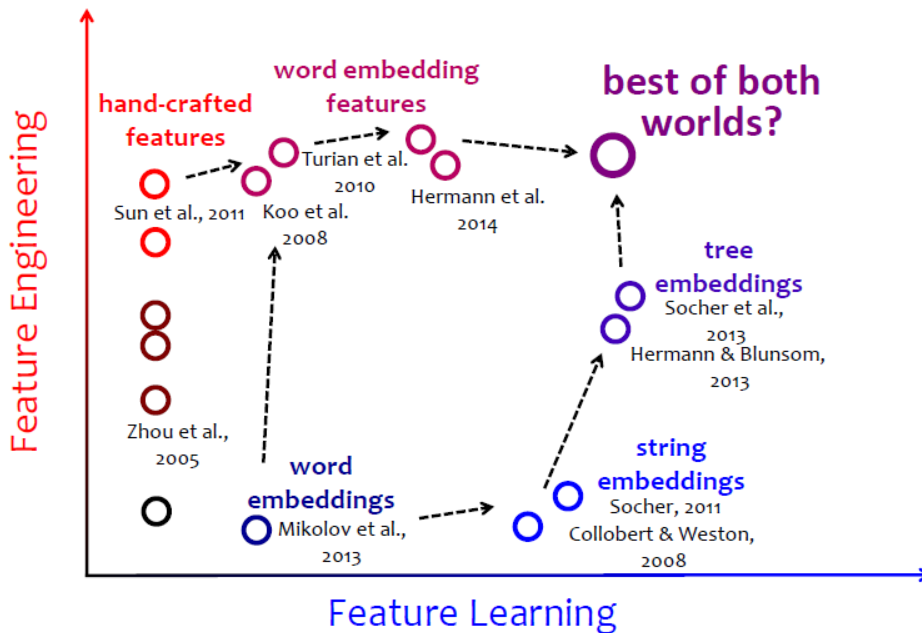
Where do features come from?



في حد فكر في انو يستخدم ال word embeddings .. ايه يعني لو استخدمنا الماشين ليرننج انو يتعلم ال vector representations of the words .. الفكرة ديه كانت جامده جداً وبقت تطير .. الناس بقت بتقول ان القطه ممكن تبقا فيكتور .. و الكلب ممكن بيقا فيكتور .. وعشان دول كلمتين similar to each other .. الفيكتر بتاعهم هيبقا شبه قريب من بعض .. المهم حد بعدها الدنيا شافها بتجري بسرعه اوي قال وليه مجريش انا كمان .. يلا بينا ندعك في اللي كلهم بيدعكو فيه .. وطلع وقال لما احنا هنعمل word embedding .. ليه منعلمش embedding to the whole sentence .. و من هنا نقدر نجيب vector representation للجمله كلها بال CNNs او بال autoencoders .. بس اصبر .. في حد تالت قال .. ايه ده احنا اصلاً عندنا ال syntax trees مرميين في كل حته من الناس اللي عملو ال feature handcraft .. ما يمكن نقدر نستخدمهم ونعمل embedding لشجره كامله في فيكتور ..

وبعدين الناس اللي كانت بتعمل handcrafting قالو .. الناس ديه شكلهم هيوصلو لحاجه ولا ايه .. تعال بسرعه ناخذ منهم الافكار ديه و نلفها ونحطها في ال
handcrafted features بتاعتنا نخليها أحسن وأحسن .. وفي الآخر الافكار بتتلاقى ..

Where do features come from?



16

تعال بقا نشوف ازاى نقدر نعمل feature handcrafting .. نفترض اننا هنبني logistic regression model .. وانت عاوز ت predict جزء من ال
speech tag لكل كلمه في الجمله .. هتستخدم انهي فينتشرز ..

Feature Engineering for NLP

Suppose you build a logistic regression model to predict a part-of-speech (POS) tag for each word in a sentence.

What features should you use?

deter. noun noun verb verb noun

The movie I watched depicted hope

17

تعال نسأل شوية أسئله عن الكلمات اللي عندنا .. وناخذها ترو وفولس مثلاً يعني

Feature Engineering for NLP

Per-word Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
is-capital(w_i)	1	0	1	0	0	0
endswith(w_i , "e")	1	1	0	0	0	1
endswith(w_i , "d")	0	0	0	1	1	0
endswith(w_i , "ed")	0	0	0	1	1	0
w_i == "aardvark"	0	0	0	0	0	0
w_i == "hope"	0	0	0	0	0	1
...

deter. noun noun verb verb noun
The movie I watched depicted hope

18

Feature Engineering for NLP

Context Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
...
w_i == "watched"	0	0	0	1	0	0
w_{i+1} == "watched"	0	0	1	0	0	0
w_{i-1} == "watched"	0	0	0	0	1	0
w_{i+2} == "watched"	0	1	0	0	0	0
w_{i-2} == "watched"	0	0	0	0	0	1
...

deter. noun noun verb verb noun
The movie I watched depicted hope

19

Feature Engineering for NLP

Context Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
...
w_i == "I"	0	0	1	0	0	0
w_{i+1} == "I"	0	1	0	0	0	0
w_{i-1} == "I"	0	0	0	1	0	0
w_{i+2} == "I"	1	0	0	0	0	0
w_{i-2} == "I"	0	0	0	0	1	0
...

deter. noun noun verb verb noun
The movie I watched depicted hope

20

Feature Engineering for NLP

Table 3. Tagging accuracies with different feature templates and other changes on the WSJ 19-21 development set.

Model	Feature Templates	# Feats	Sent. Acc.	Token Acc.	Unk. Acc.
3gramMemm	See text	248,798	52.07%	96.92%	88.99%
naacl 2003	See text and [1]	460,552	55.31%	97.15%	88.61%
Replication	See text and [1]	460,551	55.62%	97.18%	88.92%
Replication	+ rareFeatureThresh = 5	482,364	55.67%	97.19%	88.96%
5w	+ $\square_0, W_{-2} \square_0, W_2 \square$	730,178	56.23%	97.20%	89.03%
5wShapes	+ $\square_0, s_{-1} \square_0, s_0 \square_0, s_{+1} \square$	731,661	56.52%	97.25%	89.81%
5wShapesDS	+ distributional similarity	737,955	56.79%	97.28%	90.46%

adding more features

1% improvement only

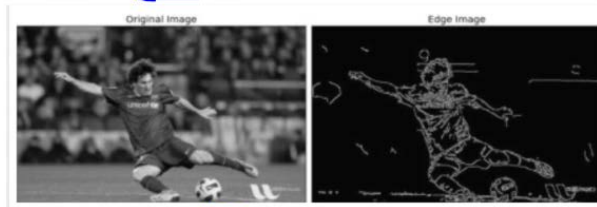
deter. noun noun verb verb noun
The movie I watched depicted hope

21

تعال نشوف شوية الكمبيوتر فيجن ..

Feature Engineering for CV

Edge detection (Canny)



Corner Detection (Harris)



ures from <http://opencv.org>

26

اللي عاوزين نتكلم عنو هو فكره عامه بحت .. اللي هو اننا نعمل some non-linear representation of your original representation .. عشان تفهم ده تعال نشوف أمثله ..

Nonlinear Features

- aka. “nonlinear basis functions”
- So far, input was always $\mathbf{x} = [x_1, \dots, x_M]$
- **Key Idea:** let input be some function of \mathbf{x}
 - original input: $\mathbf{x} \in \mathbb{R}^M$ where $M' > M$ (usually)
 - new input: $\mathbf{x}' \in \mathbb{R}^{M'}$
 - define $\mathbf{x}' = b(\mathbf{x}) = [b_1(\mathbf{x}), b_2(\mathbf{x}), \dots, b_{M'}(\mathbf{x})]$ where $b_i: \mathbb{R}^M \rightarrow \mathbb{R}$ is any function

For a linear model:
still a linear function of $b(\mathbf{x})$ even though a nonlinear function of \mathbf{x}

Examples:

- Perceptron
- Linear regression
- Logistic regression

- **Examples:** ($M = 1$)

polynomial

$$b_j(x) = x^j \quad \forall j \in \{1, \dots, J\}$$

→ radial basis function

$$b_j(x) = \exp\left(\frac{-(x - \mu_j)^2}{2\sigma_j^2}\right)$$

sigmoid

$$b_j(x) = \frac{1}{1 + \exp(-\omega_j x)}$$

log

$$b_j(x) = \log(x)$$

30

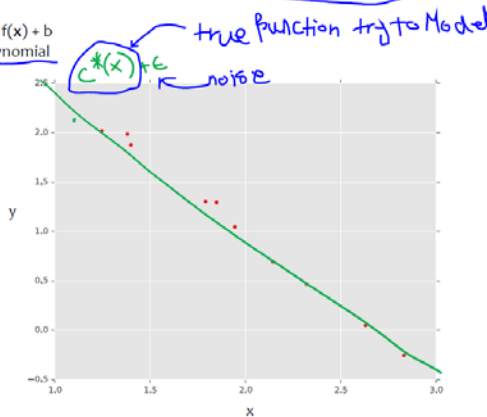
الفكره انك لو بتبني linear model .. وانت الموديل اللي انت بتتعلمو لما هتديلو $b(x)$ هو لسه linear function of $b(x)$.. احنا مش بنغير الأليجورنزم ولا الموديلز .. احنا بس عندنا non-linear function من الأنبيوت إكس بتاعنا .. فالكلاسيفير بتاعنا بقا هو فعلياً non-linear in x حتي لو كان لسه linear in $b(x)$.. تعال نشوف الكلام ده .. حتي لو هنروح لل high dimensional feature space انا اقدر لسه اوريك الفانكشن في ال low-dimensional space

Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$ where $\mathbf{f}(\cdot)$ is a polynomial basis function

y	x
2.0	1.2
1.3	1.7
0.1	2.7
1.1	1.9

true “unknown” target function is linear with negative slope and gaussian noise



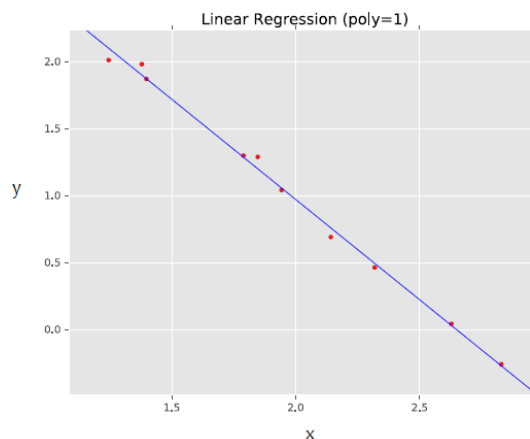
31

Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$ where $\mathbf{f}(\cdot)$ is a polynomial basis function

y	x
2.0	1.2
1.3	1.7
0.1	2.7
1.1	1.9

true “unknown” target function is linear with negative slope and gaussian noise



32

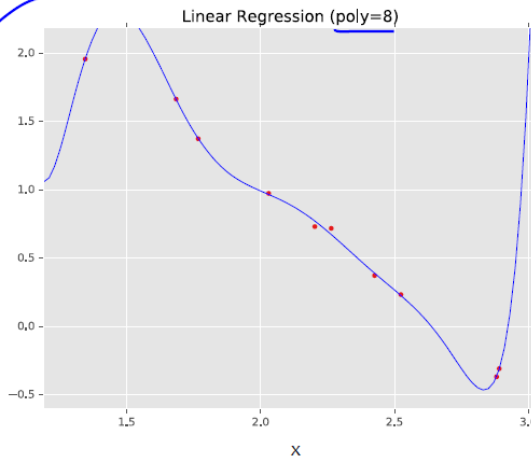
Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function

even more features

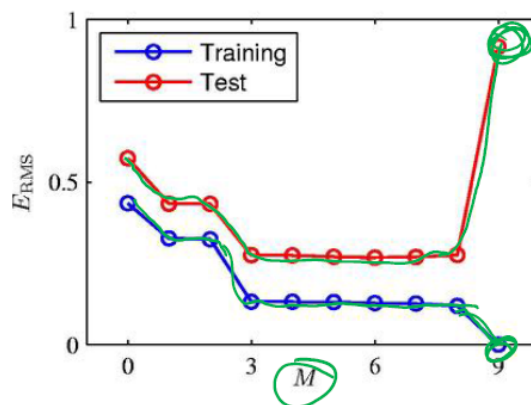
y	x	x^2	...	x^8
2.0	1.2	$(1.2)^2$...	$(1.2)^8$
1.3	1.7	$(1.7)^2$...	$(1.7)^8$
0.1	2.7	$(2.7)^2$...	$(2.7)^8$
1.1	1.9	$(1.9)^2$...	$(1.9)^8$

true "unknown"
target function is
linear with
negative slope
and gaussian
noise



36

Over-fitting



Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
θ_0	0.19	0.82	0.31	0.35
θ_1		-1.27	7.99	232.37
θ_2			-25.43	-5321.83
θ_3			17.37	48568.31
θ_4				-231639.30
θ_5				640042.26
θ_6				-1061800.52
θ_7				1042400.18
θ_8				-557682.99
θ_9				125201.43

explosion

عند $M=9$.. البامترز ظارت فوق جداً و الموديل شغال علي آخرو عشان يعمل أوفر فيتنج .. ازاي نحل المشكله ديه .. الدكتور كان مدينا 10 نقط بس .. طب لو خدنا نفس الفيتشر سيت اللي عندنا بس قمت مدليك 100 نقطه .. المشكله راحت خلاص .. ومعظم الأوفر فيتنج اتشال .. فبالتالي كل ما تقدر تصيف داتا انت هتمنع الأوفر فيتنج .. المشكله ان ايه اللي هيحصل لو معدكش داتا ... خش بالريجيو لاريزيشن ...

تعال نفكر ايه هو الأوفر فيتنج ... الأوفر فيتنج هو انو ي capture the noise in the data بدل ما ي capture the underlying structure .. ده بيحصل في كل الموديلز حرفياً .. عادي ... تعال نفكر ليه بيحصل .. ناخد مثال stock market prediction ... ايه هي الفيتشرز اللي نقدر نستخدمها .. انت تقدر تاخذ ال stock prices بتاعت كل الستوكس عند كل الاوقات .. وبعدين تبص علي كل newspaper ... و السوشال ميديا بيقلو ايه .. هل كل ده مهم؟ .. الستوكس بتاعت شركات صغيره .. ممكن كونو فعلاً بيزيدو بس هو حرفياً مفيش علاقه ما بين الي هم بيعملوه واللي جوجل بيعملوه ... فانت بتفضل دائماً ال simple hypothesis .. لما كنا بنتكلم علي ال decision trees .. احنا قلنا ان شجره صغيره هي اسهل حاجه ... تعال نقول حاجه مختلفه سيكه .. شجره صغيره ليها فيتشرز صغيرين .. فبرضو ال linear models الي بنشتغل بيهم هيقو احسن لو ليهم small number of features زي ما شفنا كذا من كام صوره فاتو ... زودت كتير اوي ع الفاضي ... في حاجه ثانيه مهمه .. وهو انك بيقا عندك important features ..

Motivation: Regularization

- **Occam's Razor:** prefer the simplest hypothesis
- What does it mean for a hypothesis (or model) to be **simple**?
 1. small number of features (**model selection**)
 2. small number of "important" features (**shrinkage**)

Handwritten diagrams illustrating model selection and shrinkage:

Setting 1: $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$, $\vec{\theta} = \begin{bmatrix} 100 \\ 0 \\ -200 \\ 0 \end{bmatrix}$

Setting 2: $\vec{x}' = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$, $\vec{\theta}' = \begin{bmatrix} 100 \\ -200 \end{bmatrix}$

Annotation: $\vec{x} \cdot \vec{\theta}$ (same product)

لو جيت تبص ع الصورة اللي فوق ديه علي طول .. هتلاقي ان ليهم نفس الدوت بروضكت ... هيبقا كويس لو فكرت ايه اللي هيحصل لو كان الصفرين اللي عندك دول هيبقوا اقرب للصفر مش صفر علي طول .. ساعتها هيبقا الدوت برضكت علي الشمال هيبقا مختلف سيكه عن اليمين .. فلما تبني موديل ليه فيتشرز كثيره .. بس معظمهم قريبين من الصفر ... انت فعلياً اتعلمت موديل سهل .. وديه الفكره ورا ال regularization .. عندك objective $J(\theta)$.. والهدف تلاقي θ اللي بتقلل فانكشن جديده .. الهدف من ال $J(\theta)$ هنا انك ت fit the data .. لو انت عملت مينيمائيزنج لل $J(\theta)$ بيقا انت اتعلمت موديل كويس ... الهدف من ال $r(\theta)$ هو بيحاول ي combat overfitting .. ازاي نختار الفورم بتاعت ال $r(\theta)$..

Regularization

- **Given** objective function: $J(\theta)$
- **Goal** is to find: $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda r(\theta)$
 (Fit the data, combat overfitting)
- **Key idea:** Define regularizer $r(\theta)$ s.t. we tradeoff between fitting the data and keeping the model simple
- **Choose form of $r(\theta)$:**

Example: q-norm (usually p-norm) $r(\theta) = \|\theta\|_q = \left[\sum_{m=1}^M \|\theta_m\|^q \right]^{\frac{1}{q}}$

	yields parameters that are...	name	optimization notes
0 $\ \theta\ _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	Lo reg.	no good computational solutions
1 $\ \theta\ _1 = \sum \theta_m $	zero values	L1 reg.	subdifferentiable
2 $(\ \theta\ _2)^2 = \sum \theta_m^2$	small values	L2 reg.	differentiable

عند $q = 0$ بتعد عدد ال non-zero entries in my theta vector . لو عملت مينيمائيزيشن لده .. فده هيفضل الارقام الكبيره خد بالك ان ده binary indicator مش سهل انك تشنقه لو انت خدت ال $q = 1$... تقدر تستخدم SGD عشان ت optimize الحوارد .. اللي مشهور أوي هو ال L2 regularizer وهو انك تاخذ ال sum of squares of the entries of the theta vector هو بيرشح ارقام صغيره في ال feature vector

Regularization

Question:

Suppose we are minimizing $J'(\theta)$ where

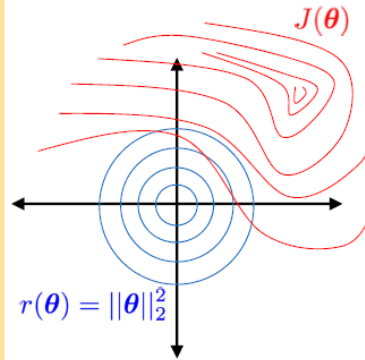
$$J'(\theta) = J(\theta) + \lambda r(\theta)$$

As λ increases, the minimum of $J'(\theta)$ will...

- A. ... move towards the midpoint between $J'(\theta)$ and $r(\theta)$
- B. ... move towards the minimum of $J(\theta)$
- C. ... move towards the minimum of $r(\theta)$
- D. ... move towards a theta vector of positive infinities
- E. ... move towards a theta vector of negative infinities

69%

~~F. ... stay the same~~ calamity



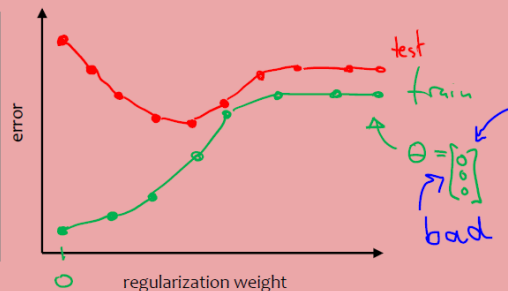
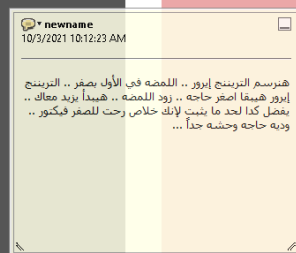
48

لو اللمضه كانت بصفر .. المنيمم هيبقا المنيمم بتاع ال $J(\theta)$ لو اللمضه كانت أكبر من الصفر سيكه ... كإنك بتجمع الاتنين مع بعض .. انت هتروح ناحية المنيمم بتاع ال $r(\theta)$ الي هو الصفر.. لما اللمضه تزيد شوية هتشد ناحية الصفر ... اللمضه تكبر أوي أوي .. ال $J(\theta)$ مقارنة بال $r(\theta)$ هتبقا صغيره جداً ... كل ما اللمضه تكبر انت رايح ناحية فيكتور من الاصفار ..

Regularization Exercise

In-class Exercise

1. Plot train error vs. regularization weight (cartoon)
2. Plot test error vs. regularization weight (cartoon)



50

بالنسبه للتست ايرور هيبداً عالي لإنك جاي من الأوفر فيتنسج وينزل معاك لإن البارمترز اللي كانت في رينج المليون عماله تقل ... فانت بتبقا احسن و احسن واحسن لحد ما تلاقي المنيمم وبعدين شوية وتثبت ...

أول حاجه don't regularize the bias parameter لازم يفكس لل bias parameter .. يعني اه انت بتضيف لل x بتاعتك 1 علي اساس انك هتضيف ال bias جوا الثيتا .. بس انت لازم متضفض عليه ترم ال regularization ... ليه ... فكر في الريجریشن موديل .. انت هيبقا عندك learning algorithm هيبقا invariant to a shift in y values ... فلو عندك شوية نقط لل regression ... فهتقدر تشيفتهم لتحت .. غير كذا هيبقو زي بعض .. السلوب نفسو انت مش عاوز يتغير انما انت بتعوز تتحرك فوق و تحت .. ده برضو بيتطبق علي الريجریشن .. انت عاوز تشيفت الديسجن باوندرى

Regularization

Don't Regularize the Bias (Intercept) Parameter!

- In our models so far, the bias / intercept parameter is usually denoted by θ_0 -- that is, the parameter for which we fixed $x_0 = 1$
- Regularizers always avoid penalizing this bias / intercept parameter
- Why? Because otherwise the learning algorithms wouldn't be invariant to a shift in the y-values ^{for regression} _{x-values for classification}

Whitening Data

- It's common to *whiten* each feature by subtracting its mean and dividing by its variance
- For regularization, this helps all the features be penalized in the same units (e.g. convert both centimeters and kilometers to z-scores)

52

انا فصلت في آخر 10 دقائق