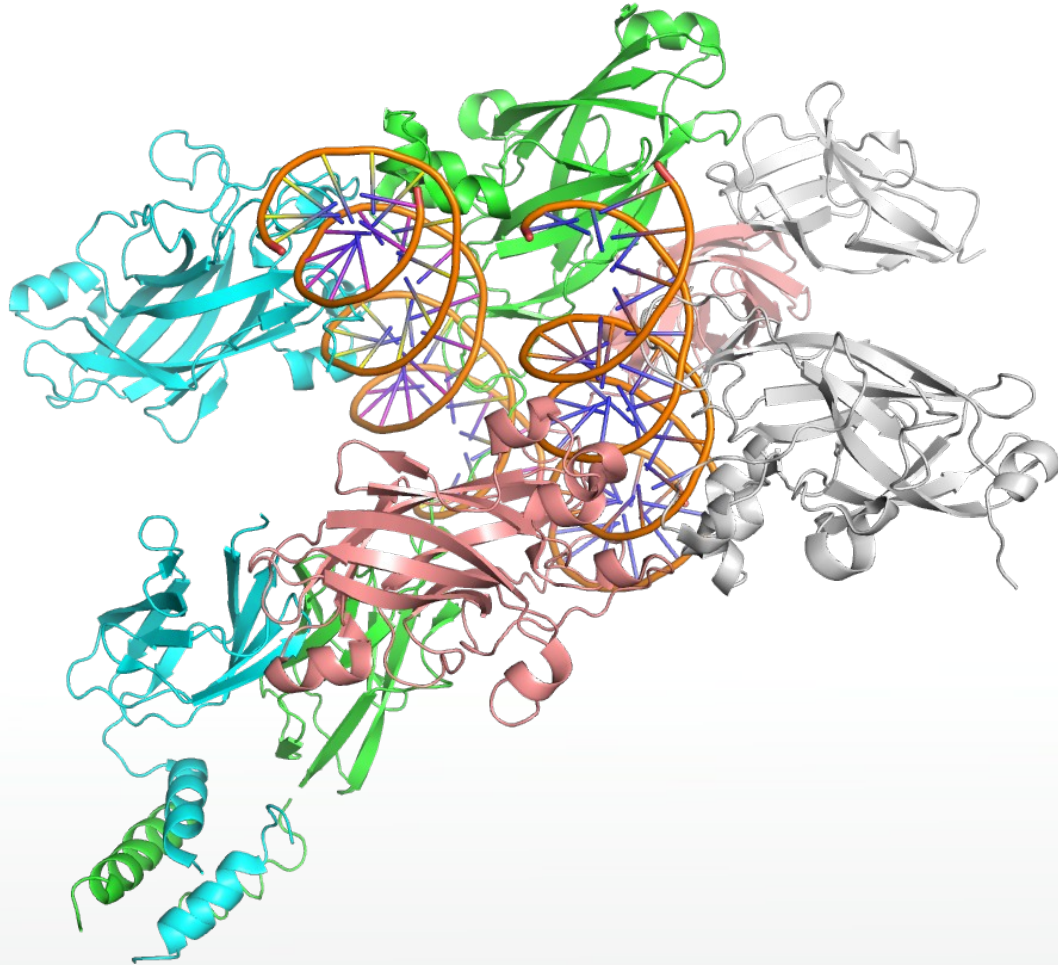




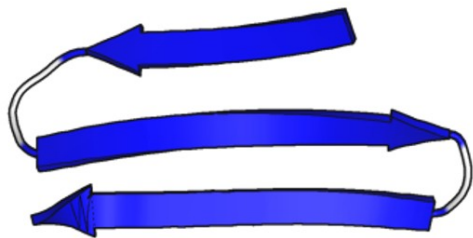
5 6 things I have learned using neural networks (deep or not)

Claudio Mirabello

What people are trying to predict about proteins



Prediction classic #1: Secondary Structure



β -Sheet (3 strands)



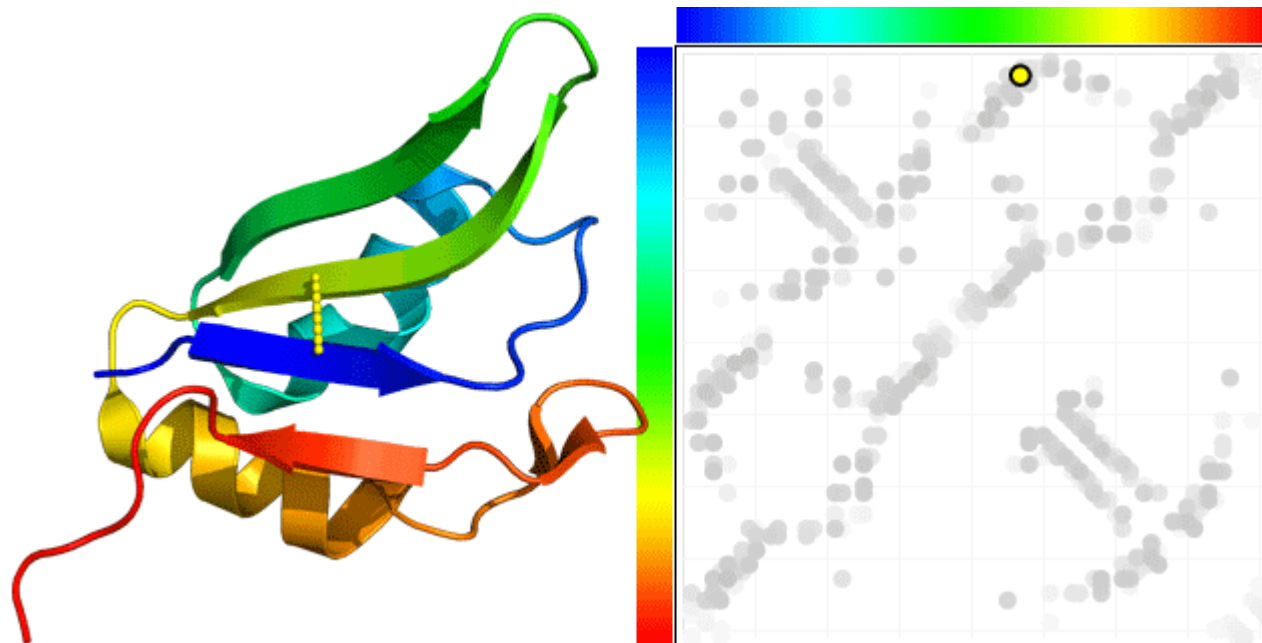
α -helix

(or none of the above)

Prediction classic #2: Solvent Accessibility

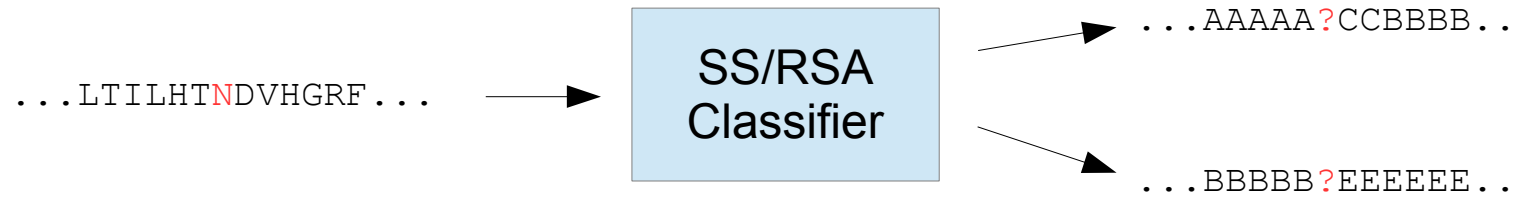


Prediction classic #3: Contact Maps



http://gremlin.bakerlab.org/gremlin_faq.php

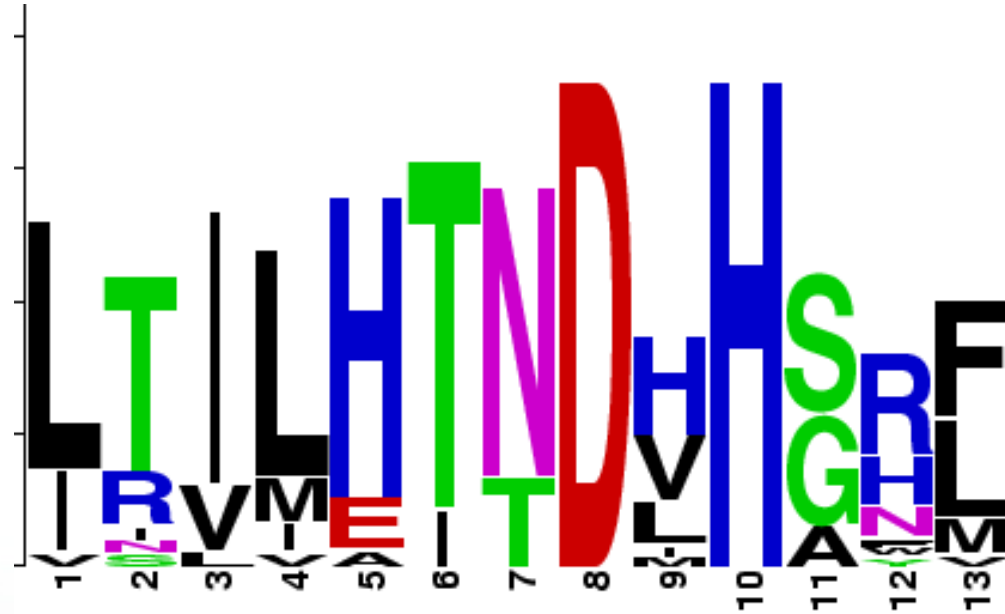
Old approach



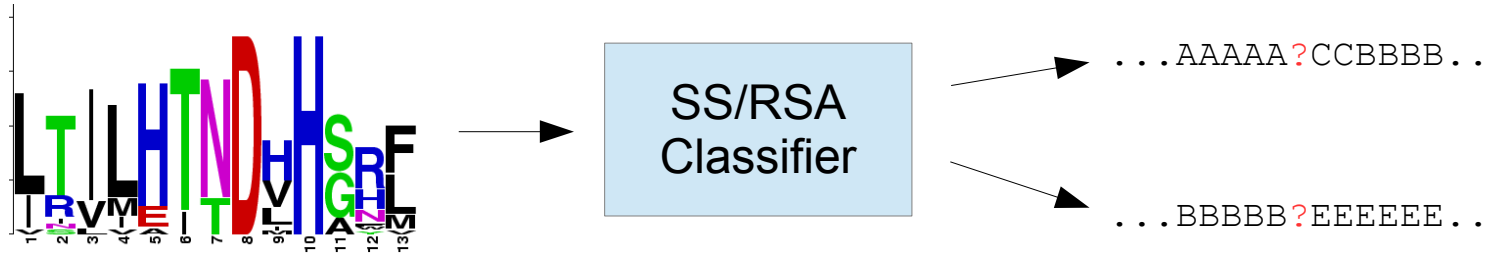
Current approach

CAV3 Hs NP 203123	--M-MAEEHTD-LEAQIVKDIHCKEIDLVRNDRKKNINEDIVKVDVFEDVIAEPVGTYSFDGVWVKVSYTTFTVSKYWCYRLI
CAV3 Mm NP 031643	--M-MTEEHTD-LEARIKDIHCKEIDLVRNDRKKNINEDIVKVDVFEDVIAEPGTYSFDGVWVKVSYTTFTVSKYWCYRLI
CAV3 Rn NP 062028	--M-MTEEHTD-LEARIKDIHCKEIDLVRNDRKKNINEDIVKVDVFEDVIAEPGTYSFDGVWVKVSYTTFTVSKYWCYRLI
CAV3 Gg NP 989701	--M-AEEQREL-EERIIKDIHCKEIDLVRNDRKKNINEDIVKVDVFEDVIAEPVGTYSFDGVWVKVSYTTFTVSKYWCYRLI
CAV3 Dr NP 991301	--M-ADQYNT--NEEKILRDSHTKEIDLVRNDRKKNINEDIVKVDVFEDVIAEPDGT-HSMDGVWVKVSYTTFTVSKYWCYRV
CAV3 Tn CAG06808	--M-ADQYQYNANEKIVKDSHTKEIDLVRNDRKKNINEDIVKVDVFEDVIAEPDGT-HSLDGVWVKVSYTTFTVSKYWCYRV
CAV3 XL AAH41289	--MAQIQQPEPAKQDKSNTKEALTEIDLVQRDEKKNINQEVVQVDFEDVIAEPDGT-HSFDGVWVKVSSSTFTVTKYWCYRV
CAV1b Dr AAN06979	--MDNDSIN----EKTLDVHTKEIDLVRNDRKKNLNDVVKVDVFEDVIAEPAGTYSFDGVWVKVSSSTFTVTKYWCYRLI
CAV1 Hs NP 001744	NNKAMADE----LSEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Bt NP 776429	NNKAMAE-----MNEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Ss NP 999603	NNKAMAE-----MNEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Oa AAT81146	NNKAMAE-----MNEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Cf NP 001003296	NNKAMAE-----MSEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Fc AAR16230	NNKAMAE-----INEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Mm NP 031642	NNKAMADE----VTEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1a Rn NP 113744	NNKAMADE----VNEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Dr NP 997816	NNKEMDND--INEKTLQDVHTKEIDLVRNDRKKNLNDVVKVDVFEDVIAEPAGTYSFDGVWVKVSSSTFTVTKYWCYRLI
CAV1 Gg AAR16241	NNKMADE----LSEKAVHDVHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Dr AAR16334	NNKEMDND-S--INEKTLQDVHTKEIDLVRNDRKKNLNDVVKVDVFEDVIAEPAGTYSFDGVWVKVSSSTFTVTKYWCYRLI
CAV1 Oc AAM19213	---AMADE----VNEKQVYDAHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPGT-HSFDGIWVKASFTTFTVTKYWFYRLI
CAV1 Tr AAR16280	NNKMDND-S--LNEKSMEDVHTKEIDLVRNDRKKNLNDVVKVDVFEDVIAEPAGTYSFDGVWVKVSSSTFTVTKYWCYRLI
CAV1 Tn AAR16324	NNKMDND-S--LNEKSLDVDHTKEIDLVRNDRKKNLNDVVKVDVFEDVIAEPAGTYSFDGVWVKVSSSTFTVTKYWCYRLI
CAV XL AAH70672	NNKTMADD-F--LTETEVRDSHTKEIDLVRNDRKKNLNDVVKIDFEDVIAEPDGT-HSFDGIWVKVSSSTFTVTKYWFYRLI
CAV2 Hs NP 001224	VQLFMDDDSYSHHSGLEYADPEKFDSDQDRDPHRLNSHL-KLGFEDVIAEPVTT-HSFDKVVWICSHALFEISKYVIMYKF
CAV2 Bt NP 001007809	VQLFMDDDSYSRHSSVDYADPKFVDFGSDRDPHRLNSHL-KVGFEDVIAEPVST-HSFDKVVWICSHALFEMSKYVIMYKF
CAV2 Mm NP 058596	VQLFMADDAYSHHSGVDYADPEKYVDSQDRDPHRLNSHL-KLGFEDVIAEPETT-HSFDKVVWICSHALFEISKYVIMYKF
CAV2 Ss AAR16299	VQLFMDDDSYSRHSGVDYADPKFVDFGSDRDPHRLNSHL-KVGFEDVIAEPVST-HSFDKVVWICSHALFEISKYVIMYKF
CAV2 Fc AAR16229	VQLFMDDDSYSRHSGVDYADPKFADSGSDRDPHRLNSHL-KVGFEDVIAEPVST-HSFDKVVWICSHALFEISKYVIMYKF
CAV2 Gg NP 001007087	TRIFMDDDNFPP--GGPALSEGEKCAEDGLERDPGLNHL--QLGFEDVIAEPILT-HSFDKVVWICSHALFELSKYVLYKL
CAV2 Dr NP 001002150	TRVIMDEDEFNRSIEPILGKKPNVYSEVQDRDPKDINKHL-KVGFEDVIAEPNST-HSFDKVVWIGSHAVFELVKYVYFRI
CAV2 Rn NP 571989	VQLFMADDAYSHHSGVDYADPEKYVDSQDRDPHRLNSHL-KLGFEDVIAEPETT-HSFDKVVWICSHALFEISKYVIMYKF
CAV2 Cf NP 001003066	VQLCMDDDAYSHHSAVDGDLQDLADSGSDRDPHRLNSHL-QVGFEDVIAADVST-HSFDKVVWICSHALFELSKYVIMYKF
CAV2 Tn AAR16323	TSIIMDEDEFNRSIEPILSQKAKASSAPDRDPHDLNAQL-KVGFEDVIAEPASA-HSFDKVVWIGSHAVFELVKFIFYRLI
CAV2 Tr CAA09081	CSIIMDEDEFNRSIEPILSKKARLYSSAPDRDPHDLNAQL-KVGFEDVIAEPASA-HSFDKVVWIGSSATFELVKFIFYRLI
CAV Ce NP 501743	WSRCQKGEQKEENIATG-----VDLVNRDANSMNNHW-QLNFEDIFGEADSQ-HSWDCVWRLNHTVFTAVRLFIFYRLI
CAV Ch CAE69013	KCKTKKCENEQKEDHIAIG-----MDFVNRDHNGLNNHW-QHNFDIFGEADSQ-HTWEFMWRLNLSVFNWVRLFVYRF
CAV Ce NP 506357	HKPPNPMEFED-IGVKNIAPVLIHKNNMDDRDEKD-SAQYLNTSFFEVEFNEPSEQYHSIACVWTLSEFKIFIVRIYSYKLI

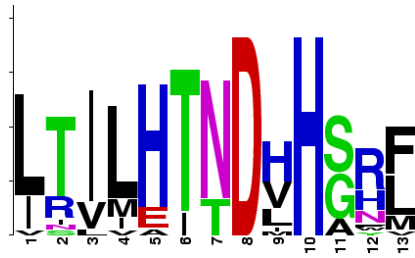
MSAs to sequence profiles



Current approach

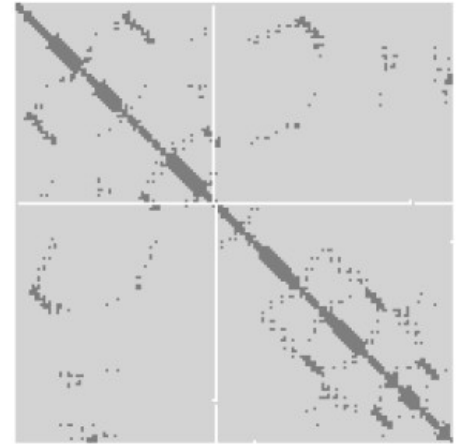


Current approach

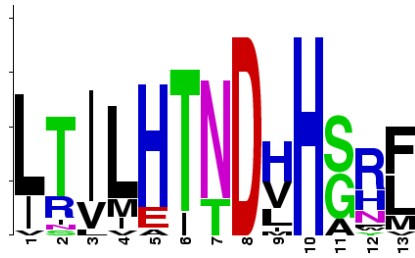


...HHHHHCCCEEEE...

Contact Map
Predictor



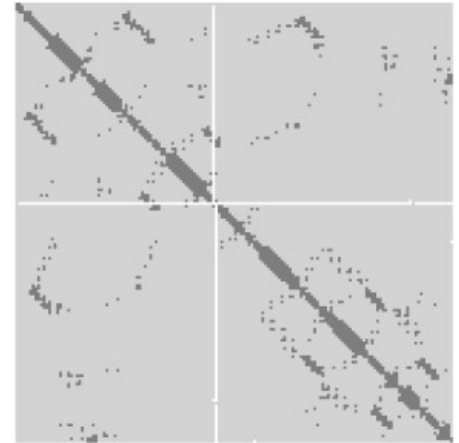
Current approach



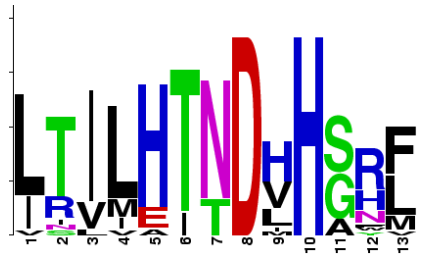
...HHHHHCCCEEEE...

...BBBBBBEEEEBB...

Contact Map
Predictor



Current approach

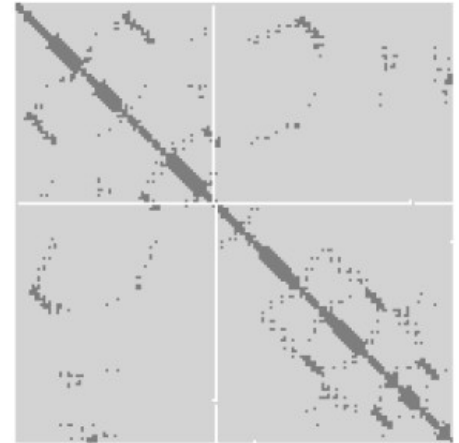


...HHHHHCCCEEEE...

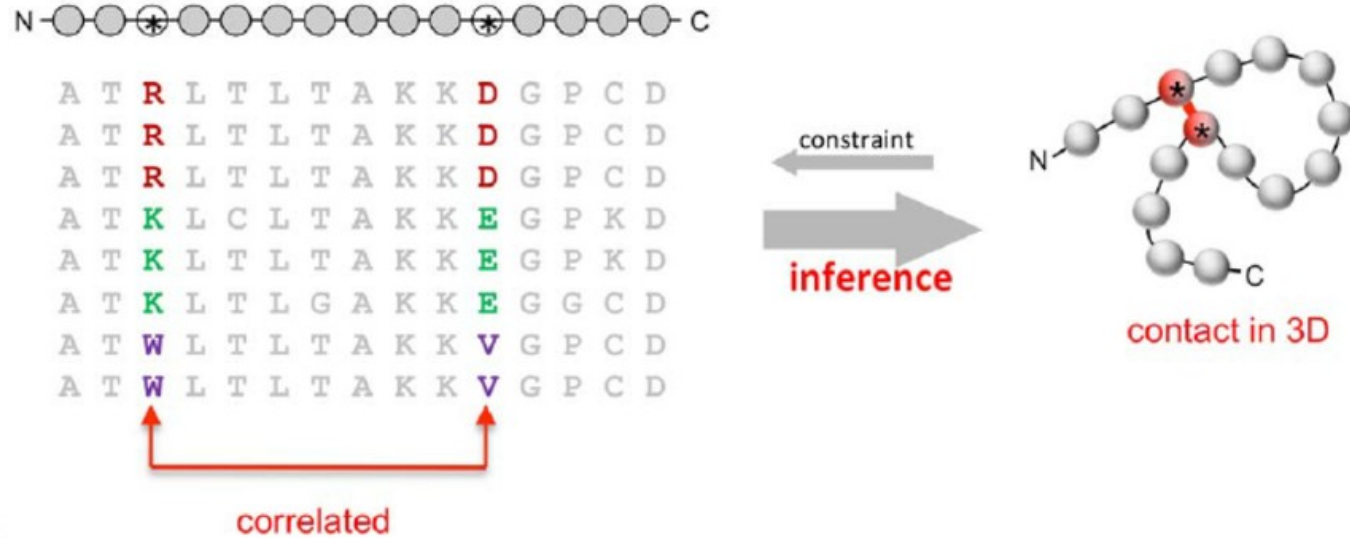
...BBBBBBEEEEBB...

Correlated mutations

Contact Map
Predictor

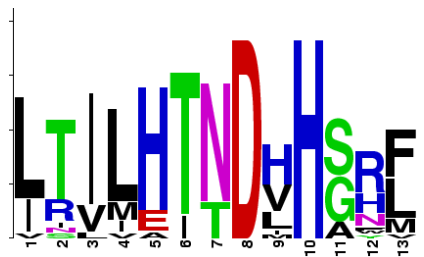


Correlated mutations



Marks, D., et al: Protein 3D Structure Computed from Evolutionary Sequence Variation (2011)

Current approach

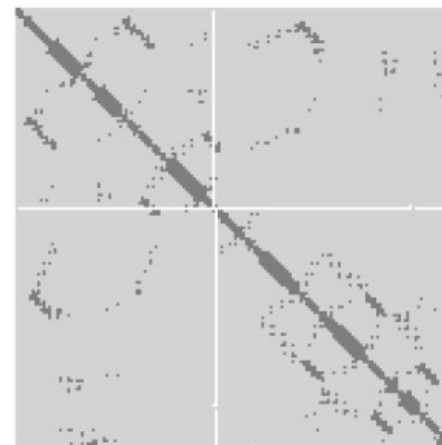


...HHHHHCCCEEEE...

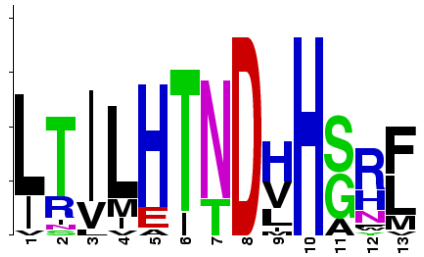
...BBBBBBEEEEBB...

Correlated mutations

Contact Map
Predictor



“Deep” approach

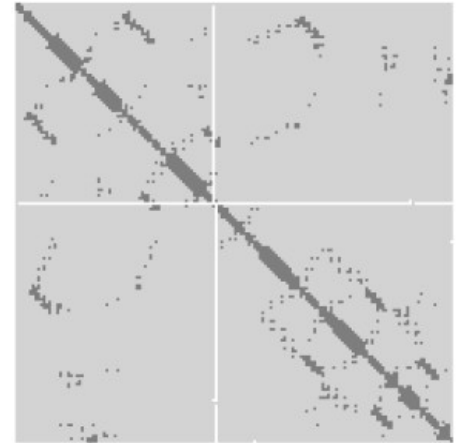


...HHHHHCCCEEEE...

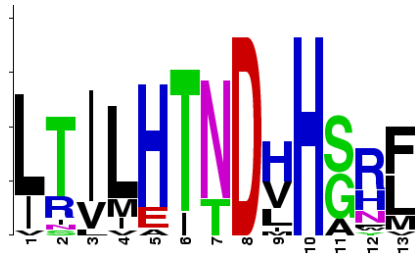
...BBBBBBEEEEBB...

Correlated mutations

Contact Map
Predictor



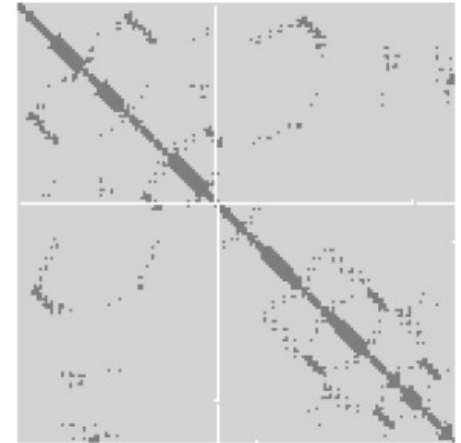
“Deep” approach



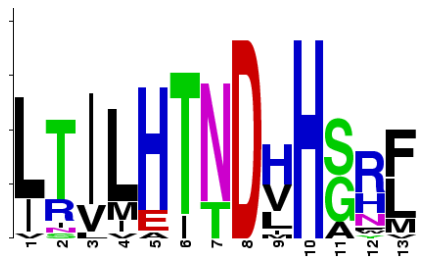
...HHHHHCCCEEEE...

...BBBBBBEEEEBB...

Correlated mutations



“Deeper” approach

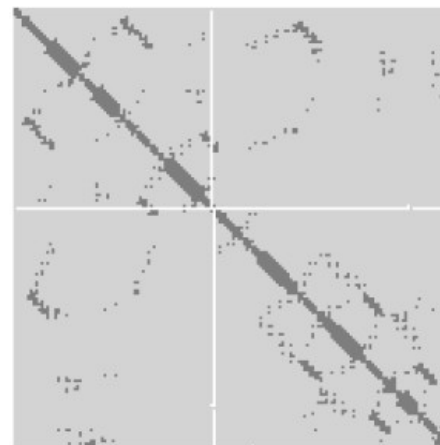


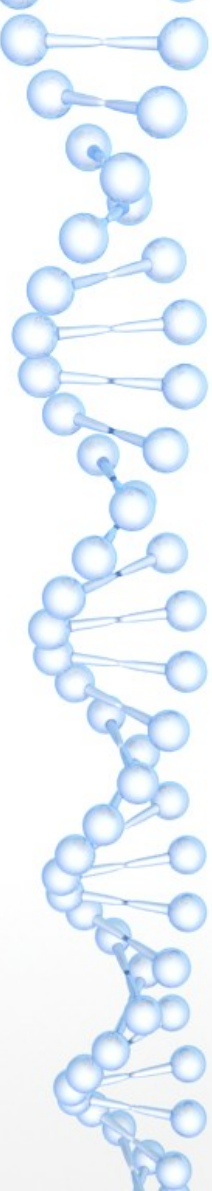
...HHHHHCCCEEEE...

...BBBBBBEEEEBB...

Correlated mutations

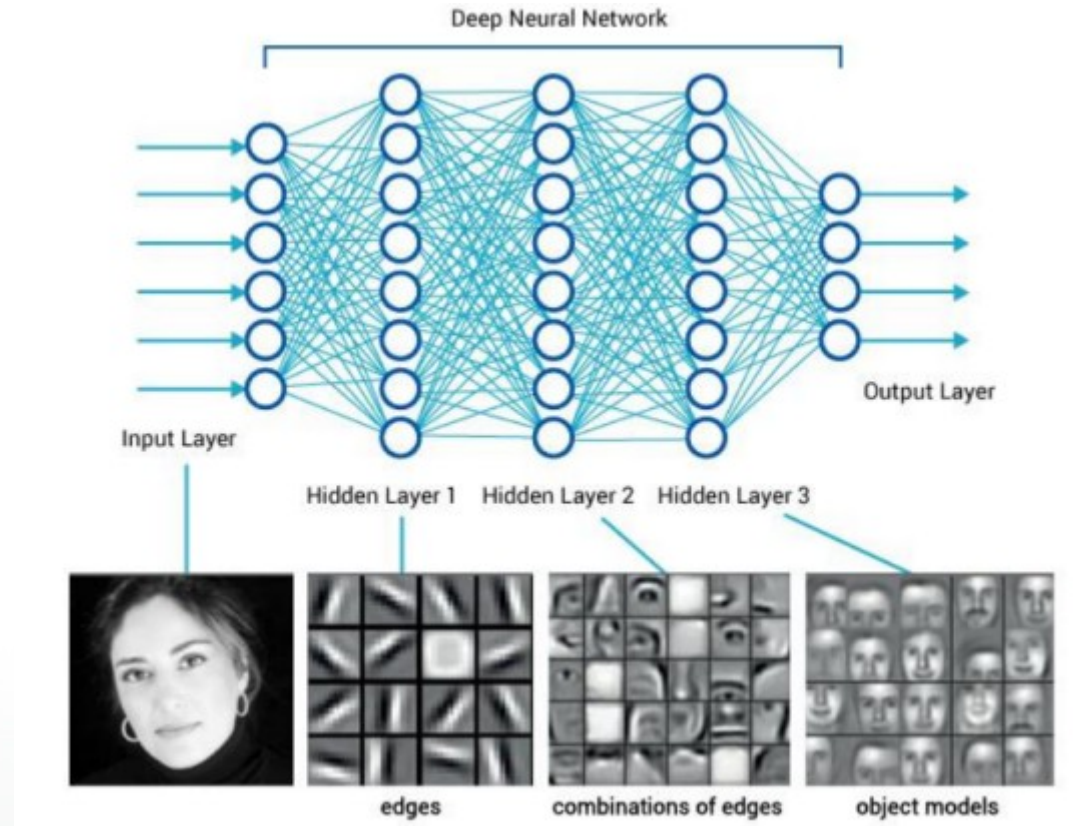
Contact Map
Predictor



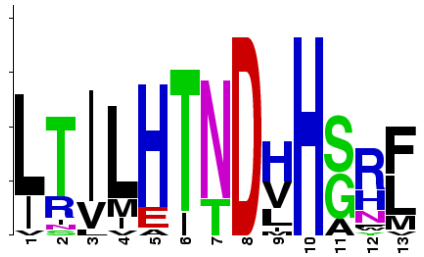


Things I've learned #1:
Deep Learning is not only about
deeper architectures,
but also about automatic feature extraction

Feature maps and automatic feature extraction



“Deep” approach

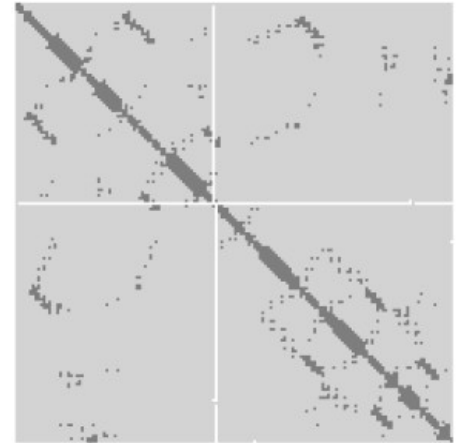


...HHHHHCCCEEEE...

...BBBBBBEEEEBB...

Correlated mutations

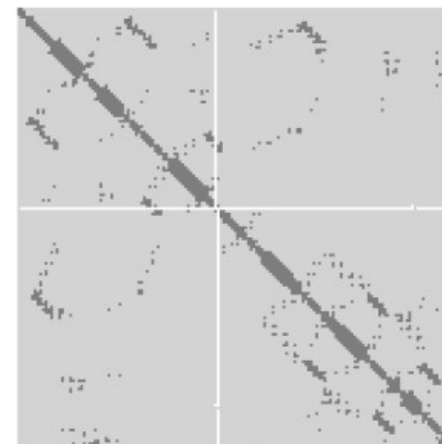
Contact Map
Predictor

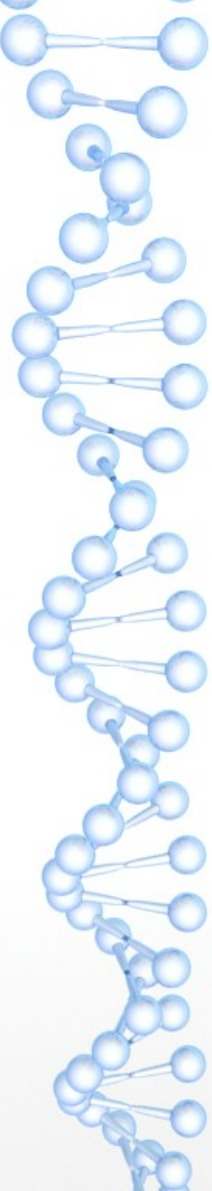


Can we input the MSA as it is?

```
CAV1 HA NP 203123 --H-HKEHTD-LEALIKDIKCKEIDLVRD KKIMEDIKVVDFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 HA NP 031643 --H-HKEHTD-LEALIKDIKCKEIDLVRD KKIMEDIKVVDFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 HA NP 062028 --H-HKEHTD-LEALIKDIKCKEIDLVRD KKIMEDIKVVDFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 GA NP 989701 --H-REYREI-REELIKDQIKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 DC NP 991301 --H-ADQVRI-REKKLEDSIKKELDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 TN CA006808 --H-ADQVRI-REKKLEDSIKKELDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 X1 AAN41289 --HAIQVRI-REKKLEDSIKKELDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV13 DQ AAN50519 --HMDRII- ----REKIQDVIKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 HA NP 001744 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 HA NP 716429 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 SA NP 999603 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 CH AAT01146 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 CE NP 001003296 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 FC AAN16236 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 HA NP 031642 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1A DQ NP 113744 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 DA NP 997616 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 GA AAN16241 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 DA AAN16334 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 CH AAN19211 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 TN AAN16280 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 TN AAN16324 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 X1 AAN70672 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV1 HA NP 001224 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 HA NP 00107809 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 HA NP 008396 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 SA AAN16299 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 FC AAN16229 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 GA NP 00100707 TRITHDDQVTF--GGLALREKCAEDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 DC NP 001002150 RVTHHDETRRLTEILAKKRVLEVDSDKIMKIL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 HA NP 571909 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 CE NP 00100366 VQLFHDDSTRISSQVADKKFVDGLDDHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 TN AAN16321 CIIHHDETRRLTEILAKKALGSADEHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 TN CA009081 CIIHHDETRRLTEILAKKALGSADEHRLMHL-KVLEFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 CH NP 001743 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 CH CA009013 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
CAV2 CH NP 506317 MKKHAD- --L-REKQVIDAKKEIDLVRD KILMIDPVKVFEDVIAE AGTT-SFDWKKVSTFTTVKVCYKRL
```

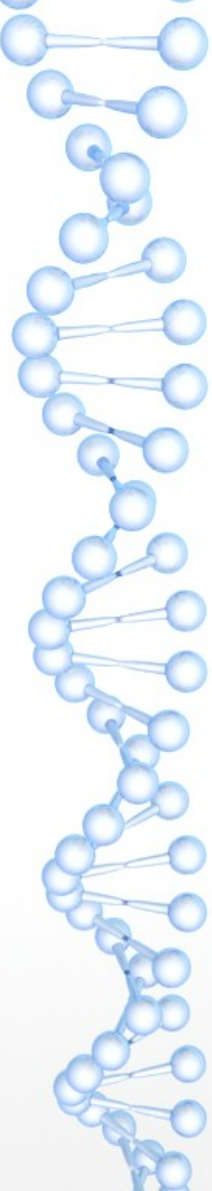
Contact Map
Predictor





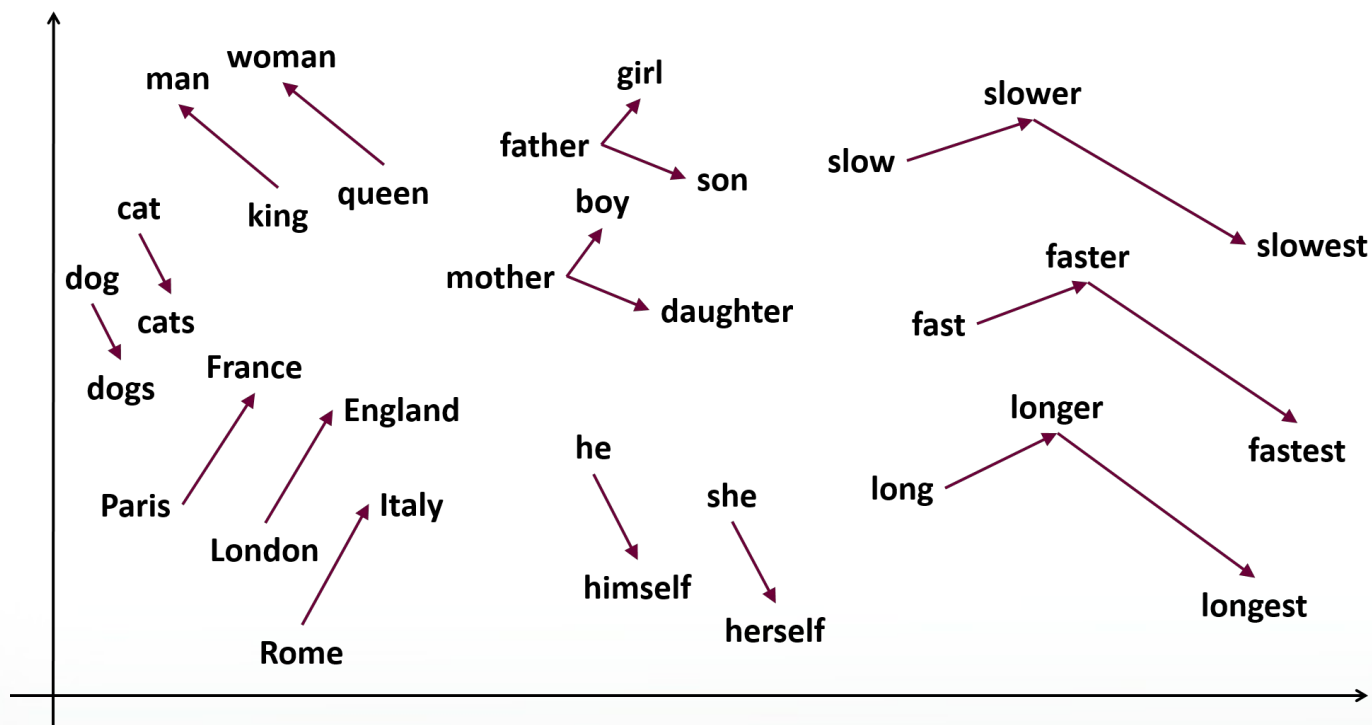
Things I've learned #2:
Don't be afraid of getting creative
when designing your model

MSAs aren't an easy input

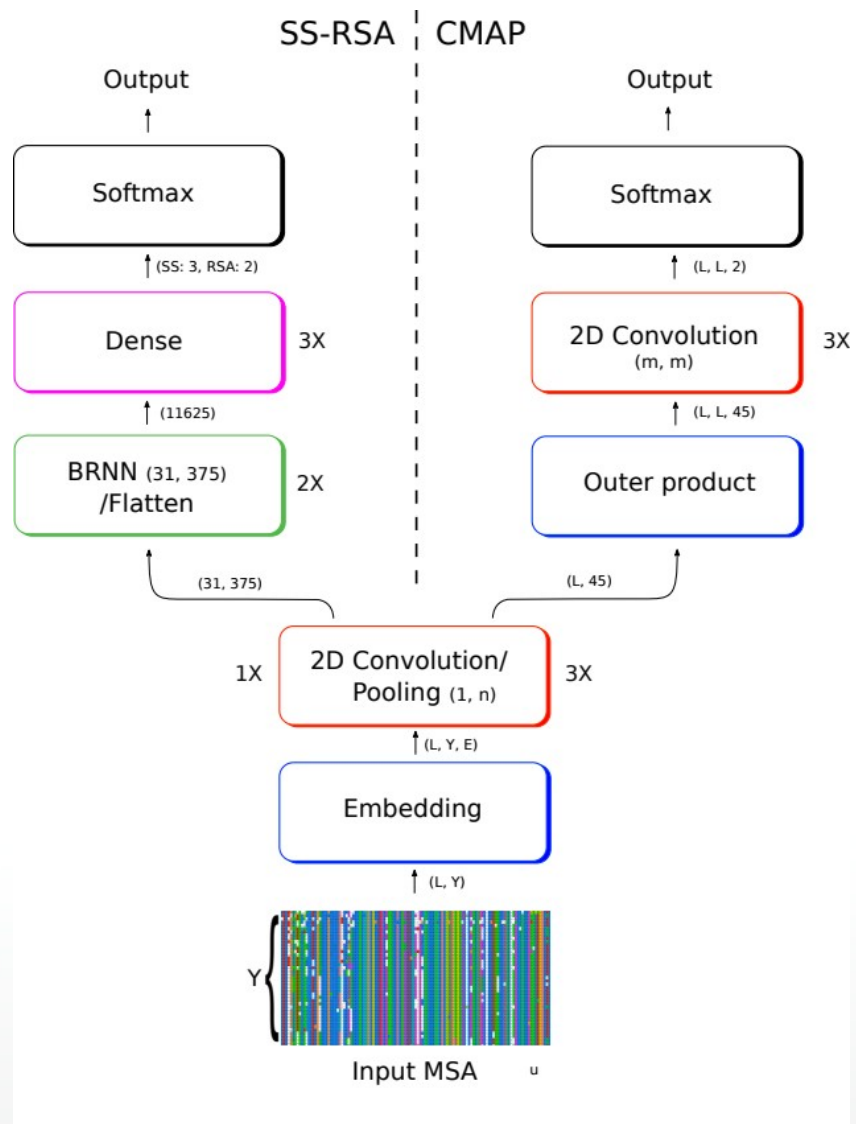
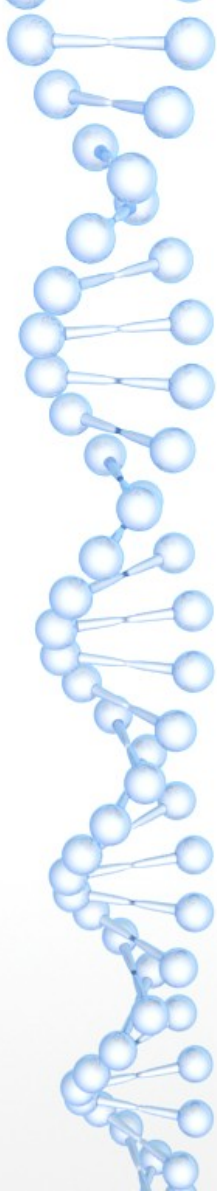


CAV3 Hs NP 203123	--M-MAEEHTD-LEAQIVKDIHCKEIDLVRNDEPKNINEDIVKVDVFEDVIAEPVGTYSFDGVWVKVSYTFTVSKYWCYRLI
CAV3 Mm NP 031643	--M-MTEETHD-LEARIIKDIHCKEIDLVRNDEPKNINEDIVKVDVFEDVIAEPVGTYSFDGVWVKVSYTFTVSKYWCYRLI
CAV3 Rn NP 062028	--M-MTEETHD-LEARIIKDIHCKEIDLVRNDEPKNINEDIVKVDVFEDVIAEPVGTYSFDGVWVRVSYTFTVSKYWCYRLI
CAV3 Gg NP 989701	--M-AEEQREL-EERIIKDKHTKEIDLVRNDEPKRINEDVVKVDVFEDVIAEPVGTYSFDGVWVKGSYTFTVSKYWCYRLI
CAV3 Dr NP 991301	--M-ADQYNT--NEEKILRDSHTKEIDLINRDEPKQINEDVVKVDVFEDVIAEPDGT-HSMDGVWVKASVYTFVSKYWCYRV
CAV3 Tn CAG06808	--M-ADQYQYNANEKIVKDSHTKEIDLINRDEPKQINEDVVKVEFEDVIAEPDGT-HSLDGVWKLSTFTVSKYWCYRV
CAV3 XL AAH41289	--MAQIQQPEPAKQDKSNTKEIDLVRNDEPKKINQEVVQVDFEDVIAEPDGT-HSFDGVWVKASSTFTVTKYWCYRV
CAV1b Dr AAN06979	--MDNDSIN----EKTLDQDVHTKEIDLVRNDEPKHLNDDVVKVDVFEDVIAEPAGTYSFDGVWVKASFTFTVTKYWCYRLI
CAV1 Hs NP 001744	NNKAMADE----LSEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Bt NP 776429	NNKAMAE-----MNEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Ss NP 999603	NNKAMAE-----MNEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Oa AAT81146	NNKAMAE-----MNEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Cf NP 001003296	NNKAMAE-----MSEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Fc AAR16230	NNKAMAE-----INEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Mm NP 031642	NNKAMADE----VTEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1a Rn NP 113744	NNKAMADE----VNEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Dr NP 997816	NNKEMDND--S--INEKTLQDVHTKEIDLVRNDEPKHLNDDVVKVDVFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Gg AAR16241	NNKEMDND--S--INEKTLQDVHTKEIDLVRNDEPKHLNDDVVKVDVFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Dr AAR16334	NNKEMDND--S--INEKTLQDVHTKEIDLVRNDEPKHLNDDVVKVDVFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Oc AAM19213	---AMADE-----VNEKQVYDAHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Tr AAR16280	NNKMDND--S--LNEKSMEDVHTKEIDLVRNDEPKHLNDDVVKVDVFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV1 Tn AAR16324	NNKMDND--S--LNEKSLQDVHTKEIDLVRNDEPKHLNDDVVKVDVFEDVIAEPVGT-HSFDGIWKASFTFTVTKYWCYRLI
CAV XL AAH70672	NNKTMADD--F--LITEVVRDSHTKEIDLVRNDEPKHLNDDVVKIDFEDVIAEPDGT-HSFDGIWKSTSTFTVTKYWCYRLI
CAV2 Hs NP 001224	VQLFMDDDSYSHHSGLEYADPEKFAADSDQDRDPHRLNSHL-KLGFEDVIAEPVTT-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Bt NP 001007809	VQLFMDDDSYSRHSSVDYADPKFVDPGSDRDPHRLNSHL-KVGFEDVIAEPVST-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Mm NP 058596	VQLFMDDDSYSHHSGVDYADPEKVFVDSHSDRDPHQLNSHL-KLGFEDLIAEPETT-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Ss AAR16299	VQLFMDDDSYSRHSGVDYADPKFVDPGSDRDPHRLNSHL-KVGFEDVIAEPVST-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Fc AAR16229	VQLFMDDDSYSRHSGVDYADPKFADSGSDRDPHRLNSHL-KVGFEDVIAEPVST-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Gg NP 001007087	TRIFMDDDNFPP--GGPALSEGEKCAEDGLERDPRGLNAHL-QLGFEDVIAEPVTT-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Dr NP 001002150	TRVIMDEDEFNRSIEPILGKKPNVYSEVQDRDPKDINKHL-KVGFEDVIAEPVST-HSFDKVVWIGSHAVFELVKYVYKFI
CAV2 Rn NP 571989	VQLFMDDDSYSHHSGVDYADPEKVFVDSHSDRDPHQLNSHL-KLGFEDLIAEPVTT-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Cf NP 001003066	VQLCMDDDSYSHHSAVDFGDLQADSGSDRDPRLNSHL-QVGFEDVIAEPVST-HSFDKVVWICSHALFEISKYVYKFI
CAV2 Tn AAR16323	TSIIMDEDEFNRSIEPILSKAKARLYSSAPDRDPHDLNAQL-KVGFEDVIAEPASA-HSFDKVVWIGSHAVFELVKYVYKFI
CAV2 Tr CAA09081	CSIIIMDEDEFNRSIEPILSKARLYSSAPDRDPHDLNAQL-KVGFEDVIAEPASA-HSFDKVVWIGSHAVFELVKYVYKFI
CAV Ce NP 501743	WSRCQKGEGEQKEENIAIG-----VDLVRNDSANMNHV-QLNFEDIFGEADSQ-HSWDCVWRLNHTVFTAVRLFYRLI
CAV Cb CAE69013	KCKTKKCENQKEDHIAIG-----MDFVNRDHNGLNHHV-QHNFDDIFGEADSQ-HTWEFMWRLNASTVFNWVRLFVYRLI
CAV Ce NP 506357	HKPPNPEMEFD-IGVKNIAPVLIHKMMMDRDRPKD-SAQYLNTSFFEVEFNPESEQYHSIACVWTLSEKIFEIVRIYSYKII

Using word embeddings to represent residues



<http://www.samyzaf.com/ML/nlp/nlp.html>



W (Tryptophan)

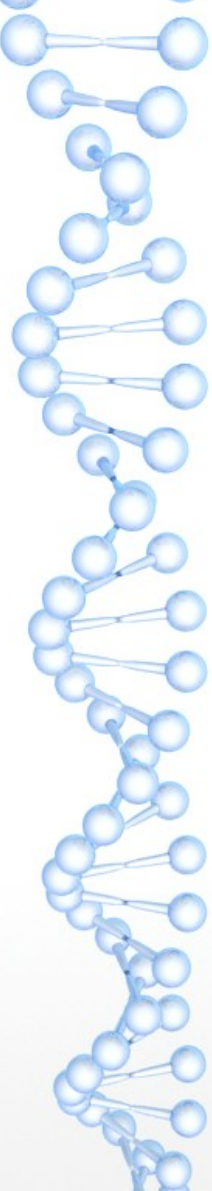
Nearest points in the original space:

I	0.072
X	0.153
Y	0.207
F	0.229
V	0.323

K

Nearest points in the original space:

H	0.029
R	0.065
D	0.182
N	0.188
.	0.201



Things I've learned #3:
Be careful with how you build your
training/testing sets

k-fold cross-validation is nice and all, but...

Training sample



Validation sample



(2F08 "Fear of Flying")



Gene

Volume 642, 5 February 2018, Pages 74-83



Research paper

Protein secondary structure prediction based on the fuzzy support vector machine with the hyperplane optimization

Shangxin Xie ^a, Zhong Li ^a , Hailong Hu ^{a, b}

 **Show more**

<https://doi.org/10.1016/j.gene.2017.11.005>

[Get rights and content](#)

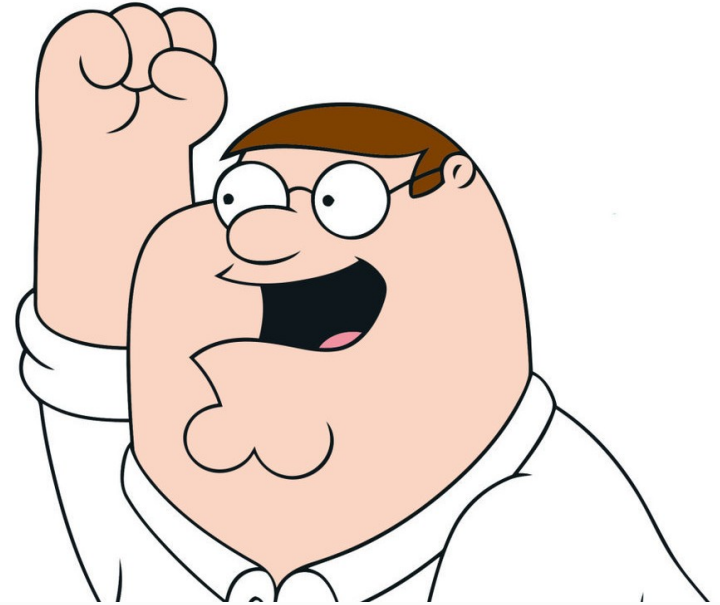
To reduce the protein sequence similarity (< 25%) between the training set and the test set, we removed 1966 protein sequences from the training set based on the sequence-based structure similarity comparison, obtaining 5986 protein sequences as the training set for the FSVM.

< 25% sequence similarity

Training sample



Validation sample





<25% sequence similarity

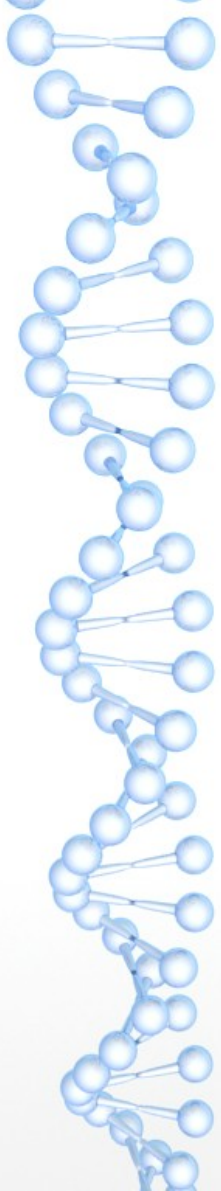
Porter, PaleAle 4.0: high-accuracy prediction of protein secondary structure and relative solvent accessibility FREE

Claudio Mirabello, Gianluca Pollastri ✉ [Author Notes](#)

Bioinformatics, Volume 29, Issue 16, 15 August 2013, Pages 2056–2058,
<https://doi.org/10.1093/bioinformatics/btt344>

Published: 14 June 2013 **Article history** ▼

resolution + $r_value/20$ or fixed at 10 for NMR structures) and redundancy, and reduced it at a **25%** sequence identity threshold, resulting in 9152 proteins. We further selected the 7522 proteins with a quality better than 4 (Full_set), but

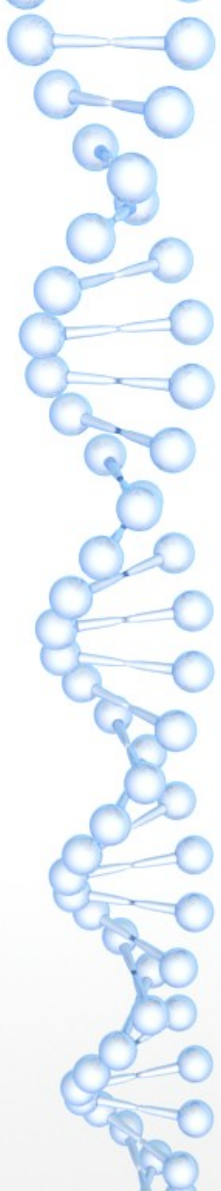


Things I've learned #4:
Not everything you try
will work for your problem
(and it will be frustrating)



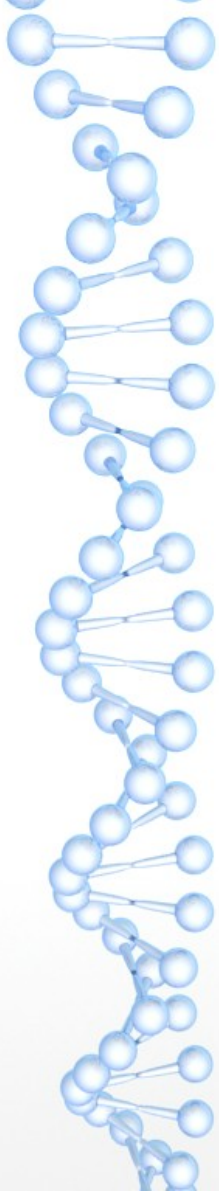
Some of the things I've tried this far

- ✓ Dropout (for SS/RSA)
- ✗ L1, L2 Regularization
- ✓ ✗ ResNets
- ✗ Batch Normalization
- ✗ Very deep networks
- ✗ Other stuff I have forgotten about

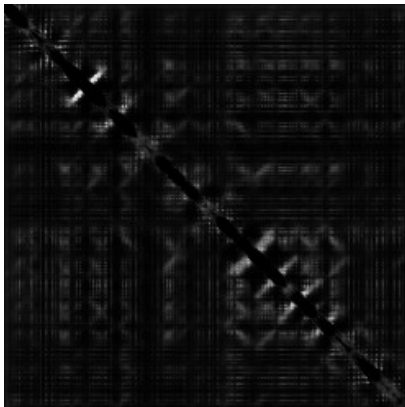


Things I've learned #5:

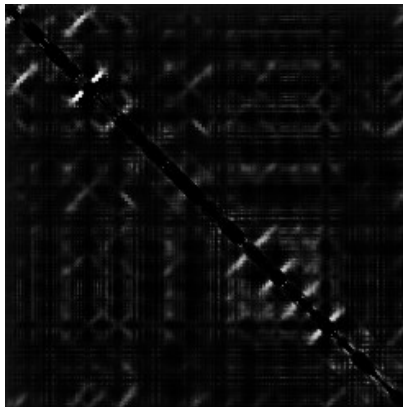
Put many weak models together and you will get a model that is better than any of them



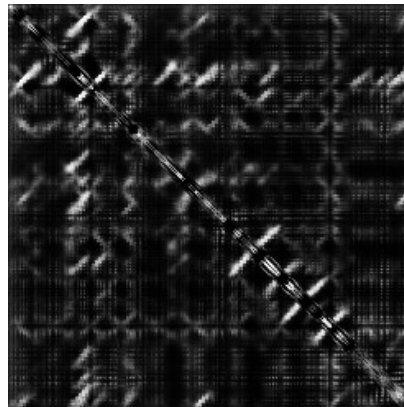
Acc: 0.271



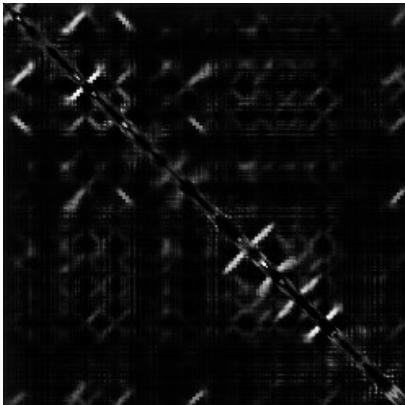
Acc: 0.277



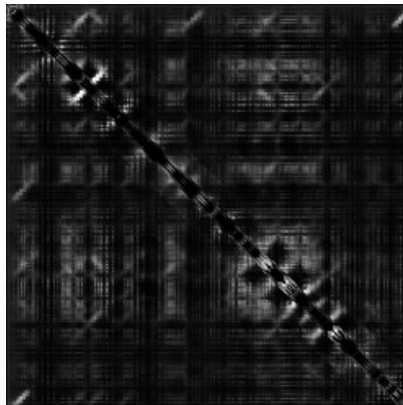
Acc: 0.346



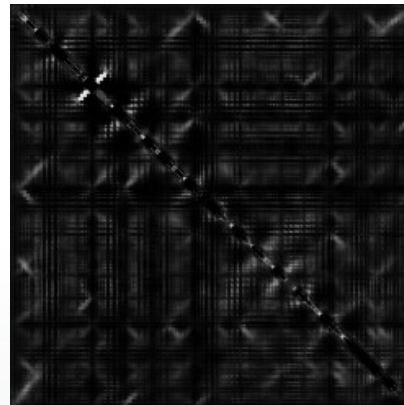
Acc: 0.340



Acc: 0.316

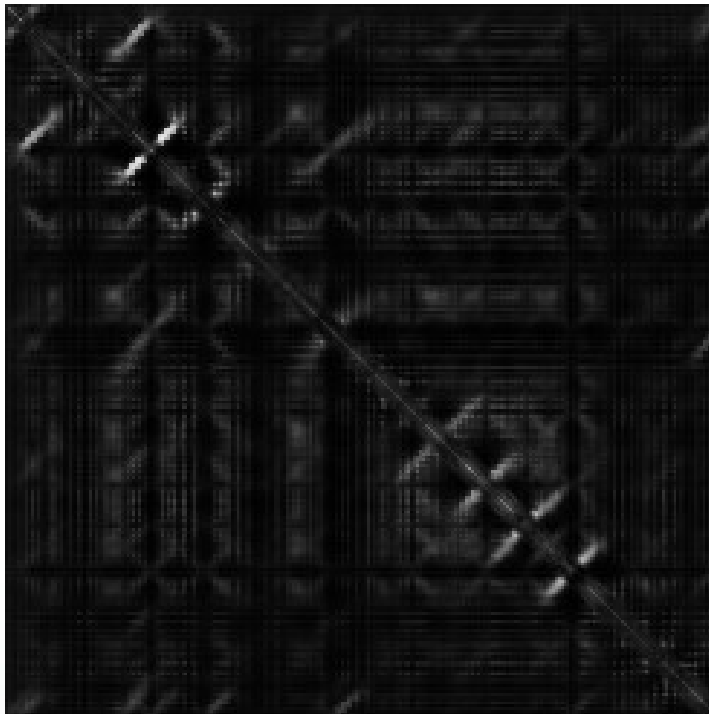


Acc: 0.289



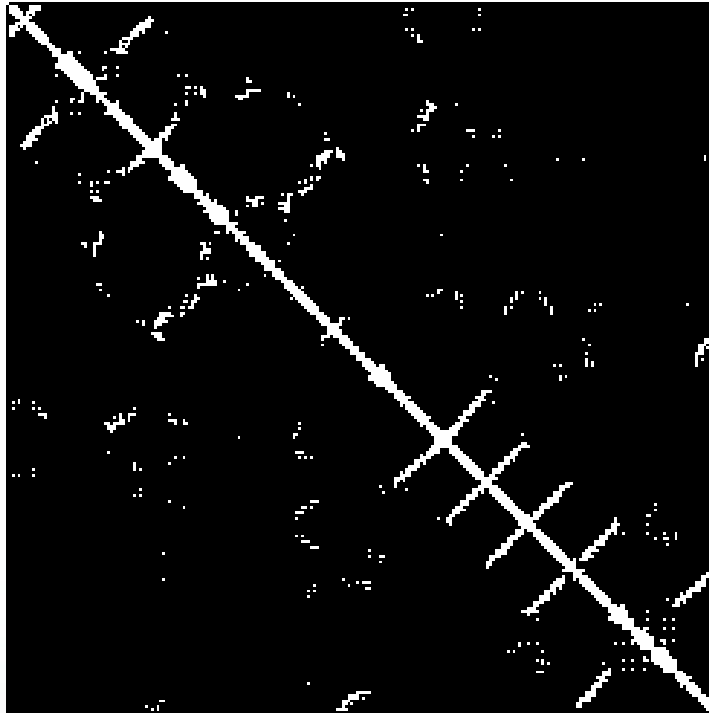
Ensembling the outputs improves the performance

Acc: 0.46

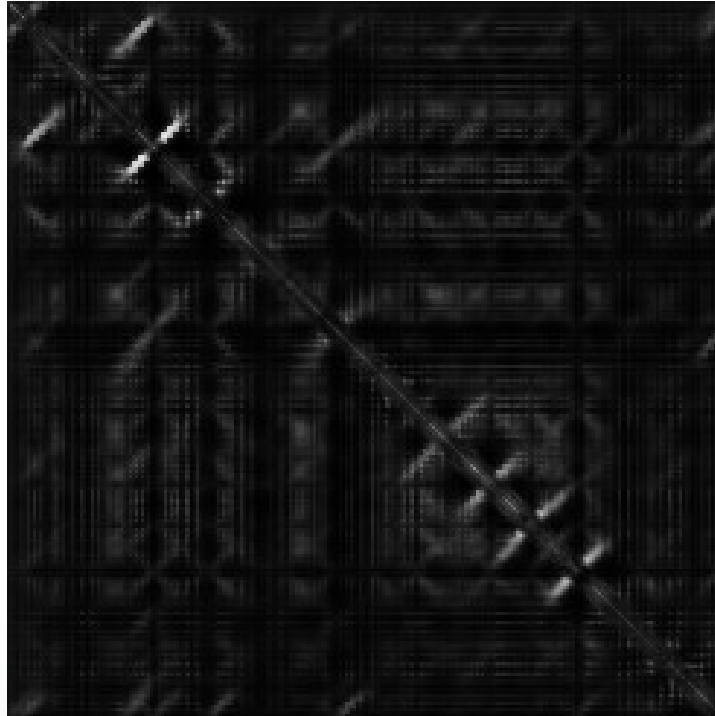


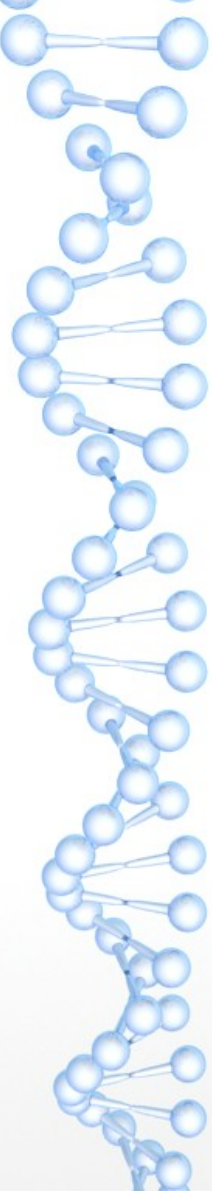
Ensembling the outputs improves the performance

Native



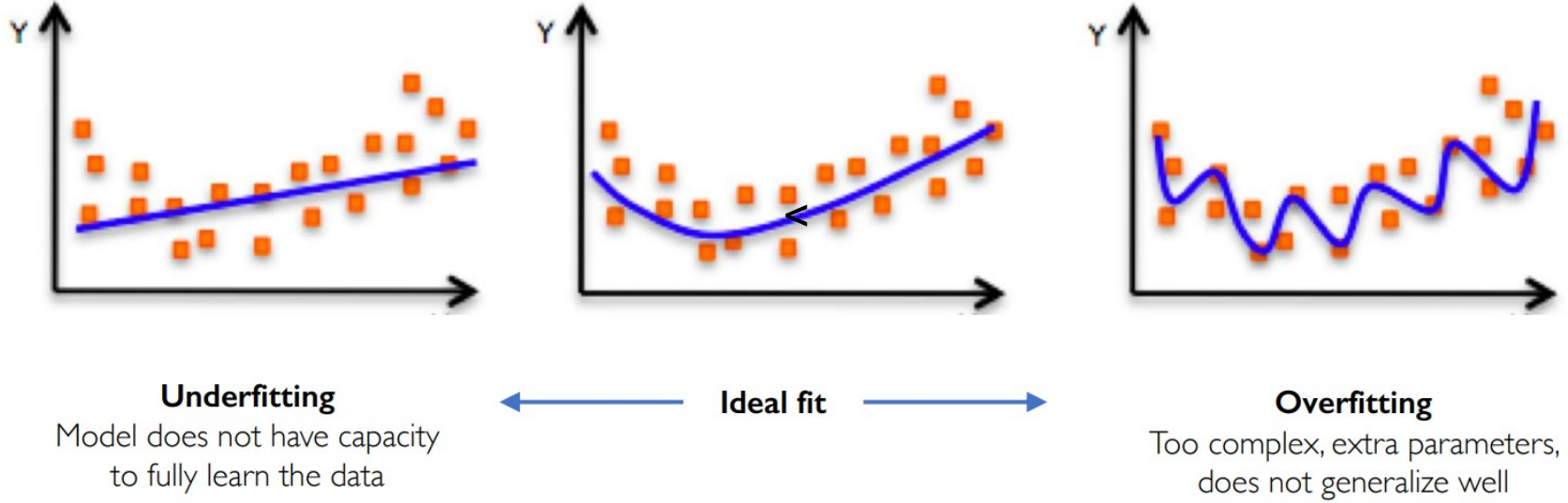
Ensembling the outputs
improves the performance



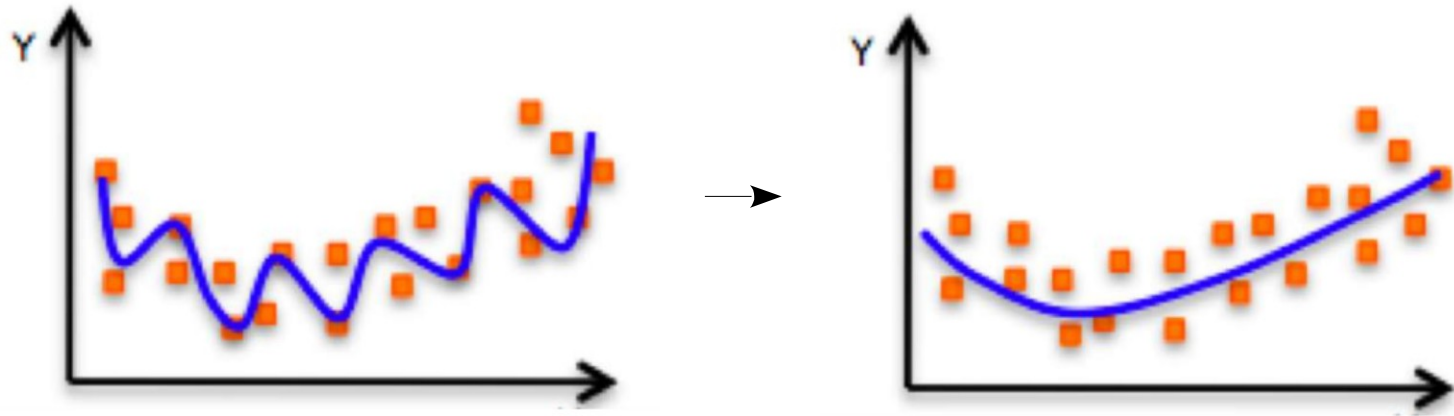


Things I've learned #6:
Overfit first, ask questions later

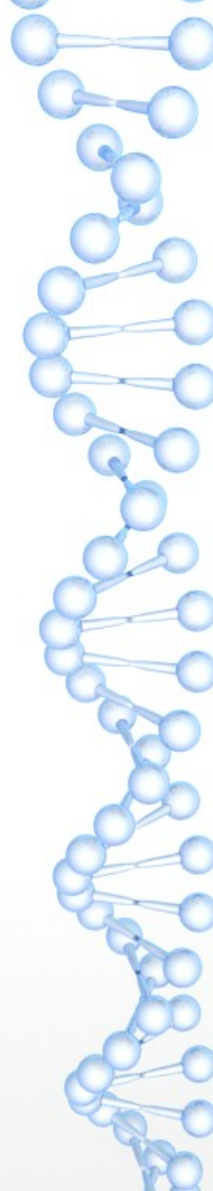
Overfitting is not necessarily a bad thing at first



If you can't overfit a model to your data
something is wrong



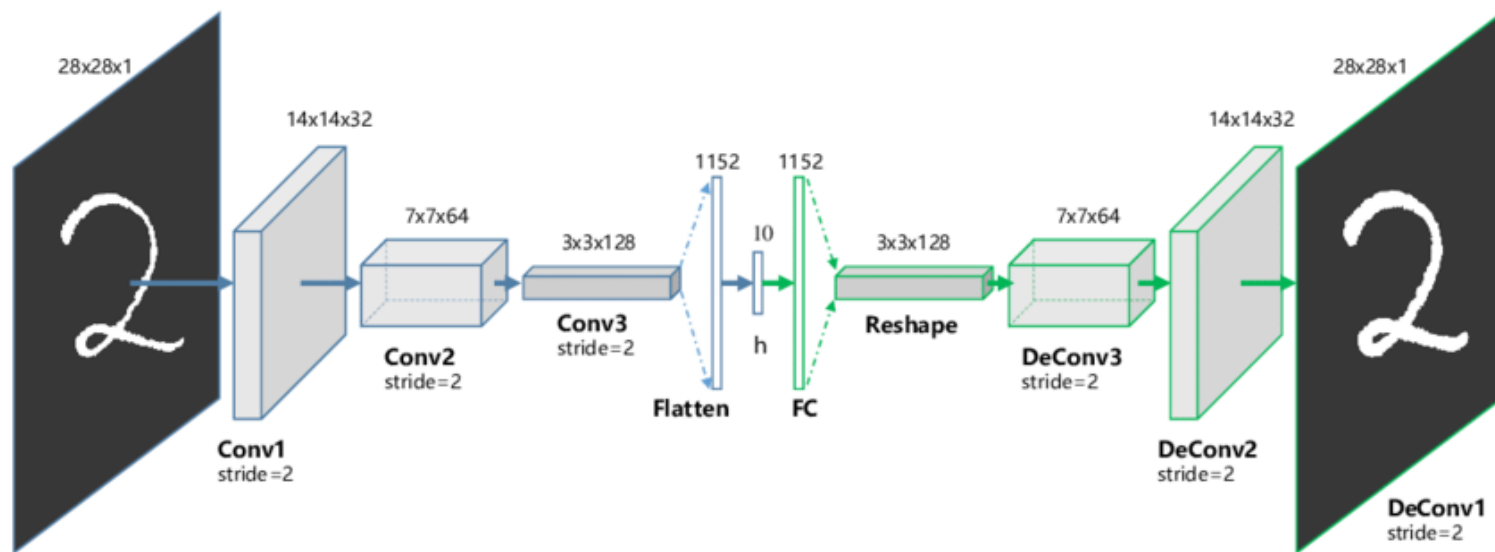
Is your model learning anything at all or making random predictions?



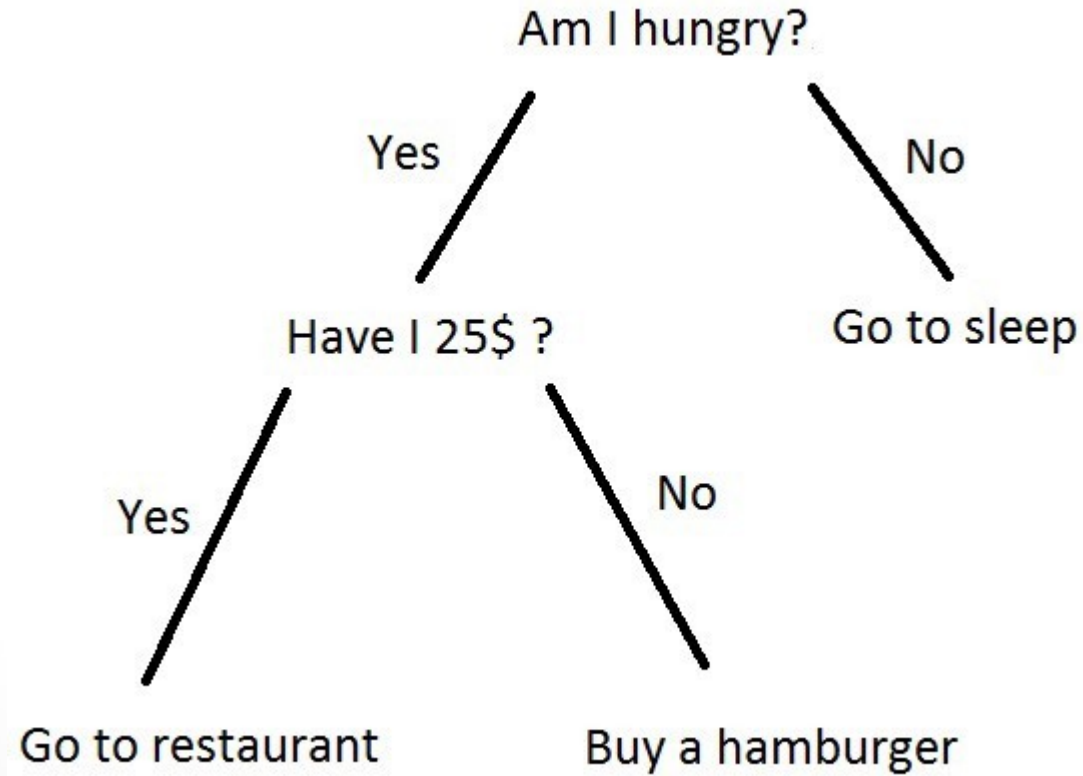
What if I have a lot of inputs and too few samples?

- ▶ Do you really need a deep network?
- ▶ Really? Are you sure?
- ▶ Rule of thumb: you want at least ~ 10 samples per each weight (Yaser S. Abu-Mostafa, Caltech)
- ▶ Reduce the number of inputs to your network, use simpler networks
- ▶ Other tricks? (more on that in a second)
- ▶ Use your domain knowledge
- ▶ Check the literature!

Autoencoders to compress inputs

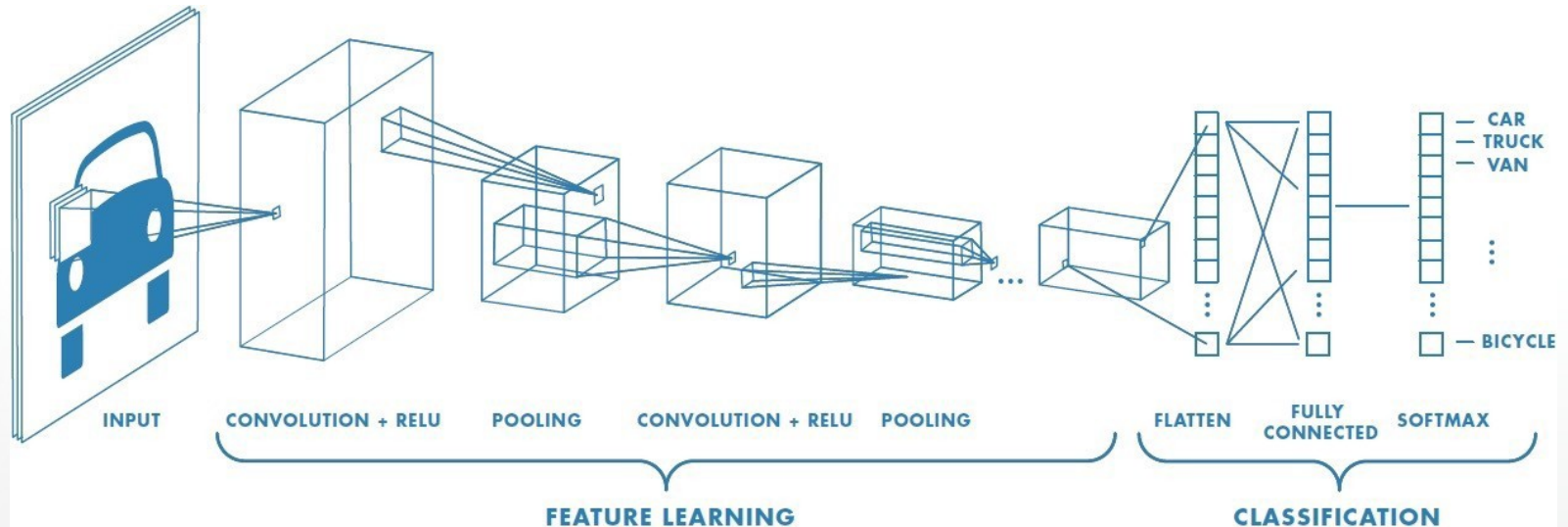


Decision trees for feature selection



Transfer learning

- ▶ Idea: use another, larger dataset containing similar data to do the initial training
- ▶ When the network has “learned” the features, you can move on to refine the training on your dataset



Transfer learning

Transfer learning: idea

