

Accelerating drug discovery with deep neural networks

literature review

Tobias Sikosek Senior Data Scientist In Silico Unit (Heidelberg)





Drug Discovery in a nutshell





Deep Learning in Drug Discovery

Learning from data to make better in silico predictions

- Target identification
 - Based on human genetic variation (DNA) associated with disease
 - Based on cellular pathways / gene expression associated with a disease
- Matching targets and small molecules with DL
 - Encode protein structure
 - Encode small molecule
 - generate new small molecules
 - Predict drug-target interactions
- Drug vs Biology: toxicity, side-effects
 - Predict toxicity of drugs from their chemical structure based on past clinical failures





protein that can be modified by drug to change disease state



Serving patient subpopulations sharing common genetic markers for disease

- Needle in a haystack problem:
 - Genome wide association studies statistically link regions within chromosomes to a particular disease / phenotype
 - Across human population, every chromosome region may contain many thousand SNVs (single nucleotide variations) which one causes the disease?
 - Often SNVs lie within DNA regions bound by transcription factors, TFs (DNA-binding proteins that act as regulatory switches within the complex circuitry that controls all cell processes)
 - If an inherited change in that DNA region leads to decreased TF binding a disease state of the cell can be the result
 - TFs are usually not direct drug targets, but may lead to the right target
- Deep Learning solution:
 - Input: DNA sequence segment
 - Output: binary classification (sequence contains TF-binding site or not)

DNA-protein binding prediction



Angermueller, C., Pärnamaa, T., Parts, L. and Stegle, O. (2016) 'Deep learning for computational biology', Molecular Systems Biology, 12(7), p. 878

Gene expression patterns reveal disease biology and pathways

- Complex network interaction problem:
 - Biology at the cellular level is the result of countless molecular interactions that can be descriped as networks (gene regulation, proteinprotein interaction, metabolic reactions, protein modifications)
 - Perturbations in this complex system (disease, environment, drugs) can have highly non-linear consequences that are difficult to model or predict
 - Cellular data contain a lot of intrinisic noise (high time-dependence, dynamics, experimental variation, etc.)
 - The most popular experimental assay to capture complex cellular biology is transcriptomics, i.e. expression (=abundance/frequency of RNA copies made from DNA gene) patterns of all ~20000 genes – or cell-type specific subset.
 - Gene expression can be highly (anti-)corellated, i.e. When high expression of a gene causes increase or decrease of a range of other genes
 - Genes can be mapped to same pathway (causally linked to a common endpoint). Example: inherited genetic change associated with a disease changes gene expression with downstream effect along the pathway. Any gene (node) in the pathway could be target of a drug intervention to modify aberrant gene expression back to normal level.

Balázsi, G., Heath, A. P., Shi, L. and Gennaro, M. L. (2008) 'The temporal response of the Mycobacterium tuberculosis gene regulatory network during growth arrest', *Molecular Systems Biology*, 4(225), pp. 1–8. ; https://commons.wikimedia.org/wiki/File:Mouse_cdna_microarray.jpg







Gene expression patterns reveal disease biology and pathways



De-noising autoencoders signal/noise from gene expression data and provide lowerdimensional fingerprint of data (\rightarrow dimensionality reduction)



Tan, J., Hammond, J. H., Hogan, D. A. and Greene, C. S. (2016) 'ADAGE-Based Integration of Publicly Available Pseudomonas aeruginosa Gene Expression Data with Denoising Autoencoders Illuminates Microbe-Host Interactions', *mSystems* 1(1), pp. e00025-15.

gsk

Gene expression patterns reveal disease biology and pathways

- Weights (parameters) between input layer (genes) and hidden layer can be used to "label" hidden nodes.
- Each hidden node is positively linked to subset of genes and negatively linked to other genes
- Each hidden node could in principle correspond to a cellular pathway (but is not restricted to any known **pathways**)
- Averaged results from **ensembles** of autoencoders yield improved results
- Outcome: which genes/pathways are most active in disease? → potential drug targets



Tan, J., Doing, G., Lewis, K. A., Price, C. E., Chen, K. M., Cady, K. C., Perchuk, B., Laub, M. T., Hogan, D. A. and Greene, C. S. (2017) 'Unsupervised Extraction of Stable Expression Signatures from Public Compendia with an Ensemble of Neural Networks', *Cell Systems*. 5(1), p. 63–71.e6.

Barcodes from L1000 gene expression (drug perturbation) - method



- L1000 data: expression of ~1000 "landmark genes" (minimal co-expression)
- Goal:
 - obtain difference profiles before and after drug treatment
 - condense information into length-100 binary B barcode
- Calculate similarity between drugs based on L1000-barcodes



Filzen, T. M., Kutchukian, P. S., Hermes, J. D., Li, J. and Tudor, M. (2017) 'Representing high throughput expression profiles via perturbation barcodes reveals compound targets', *PLOS Computational Biology*. 13(2), p. e1005335.

12

Target identification

Barcodes from L1000 gene expression (drug perturbation) - application

- New unknown compounds with verified activity against MAPK pathway were identified based on similarity of gene expression profiles to known actives
- t-SNE is a dimensionality reduction algorithm for visualization in 2D
- Z-scores are from L1000 input data
- 100D barcodes were generated by deep neural network
- Orange: known active compounds against MAPK pathway
- Circled: MAPK tool compounds









Representing drug targets at molecular detail

overview

- Most genes hold the instructions for making a particular type of protein
- Proteins are complex molecules that can be described at different levels of complexity:
 - Sequence of letters (amino acids, secondary structure)
 - List of 3D coordinates (multiple atoms per amino acid)
 - Interactions between proteins (and other molecules, e.g. drugs)



Primary structure amino acid sequence

beta sheet

alpha helix

https://en.wikipedia.org/wiki/File:Main_protein_structure_levels_en.svg; https://en.wikipedia.org/wiki/Active_site#/media/File:Enzyme_structure.svg

Encoding protein sequences



- Challenge for deep learning:
 - length of protein sequence & size of 3D structure are variable
 - machine learning models often expect fixed-length input layer
- Variable-length protein \rightarrow fixed-length input:
 - Break sequences into artificial chunks
 - Problem: often protein needs to be studied in its entirety
 - Choose input size <= longest sequence, buffer rest with "zeros"
 - Problem: wasteful

Encoding protein sequences



- ProtVec: borrows concepts from Natural Language Processing (NLP) "Word2Vec"
 - Full protein sequence ("sentence") is broken down into three-letter "words"
 - Each sentence-vector can be represented as a linear combination of word-vectors
- Treat amino acid sequence as a "sentence", AA triplets as "words"

Original Sequence (1) $\overrightarrow{M}^{(2)}\overrightarrow{A}^{(3)}\overrightarrow{F}SAEDVLKEYDRRRRMEAL..$ Splittings (1) MAF, SAE, DVL, KEY, DRR, RRM, ... (2) AFS, AED, VLK, EYD, RRR, RME, ... (3) FSA, EDV, LKE, YDR, RRR, MEA, ...

Fig 1. Protein sequence splitting. In order to prepare the training data, each protein sequence will be represented as three sequences (1, 2, 3) of 3-grams.

doi:10.1371/journal.pone.0141287.g001

Asgari, E. and Mofrad, M. R. (2015) 'Continuous Distributed Representation of Biological Sequences for Deep Proteomics and Genomics', *PLoS One*, 10(11), p. e0141287. doi: 10.1371/journal.pone.0141287.

Encoding protein sequences



- t-SNE: **2D maps** of protein space with ProtVec as input (derived from AA sequence only)
- Accurately clusters proteins based on phys-chem properties (left) and disorder (proteins with no stable structure) (right)



10.1371/journal.pone.0141287.

Predict protein structure based on sequence (and derived features)



- Input features: L=sequence length of protein
 - Sequential (L x 26)
 - Position-specific scoring matrix (20)
 - Predicted 3-state secondary structure (3)
 - Predicted 3-state solvent accessibility (3)
 - **Pairwise** (LxLx3)
 - Co-evolutionary information (CCMpred)
 - Mutual information
 - Mijazawa-Jernigan contact potential
- Output:
 - Pairwise amino-acid contact map



Wang, S., Sun, S., Li, Z., Zhang, R. and Xu, J. (2017) 'Accurate De Novo Prediction of Protein Contact Map by Ultra-Deep Learning Model', *PLOS Computational Biology*. Edited by A. Schlessinger, 13(1), p. e1005324. doi: 10.1371/journal.pcbi.1005324.

Predict protein structure based on sequence (and derived features)

gsk

Improved prediction of long-range contacts (distant along sequence, close in 3D)

Superimposition between predicted model (red) and its native structure (blue) for the CAMEO test protein (PDB ID 2nc8 and chain A).

Overlap between top L/2 predicted contacts (in red or green) and the native contactmap (in grey) for CAMEOtarget 2nc8A. Red (green) dots indicate correct (incorrect) prediction. (A) The comparison between our prediction (in upper-left triangle) and CCMpred (in lower-right triangle). (B) The comparison between our prediction (in upper-left triangle) and MetaPSICOV (in lower-right triangle).

Wang, S., Sun, S., Li, Z., Zhang, R. and Xu, J. (2017) 'Accurate De Novo Prediction of Protein Contact Map by Ultra-Deep Learning Model', *PLOS Computational Biology*. Edited by A. Schlessinger, 13(1), p. e1005324. doi: 10.1371/journal.pcbi.1005324.



Amino acid substitutions: 3D atomic coordinates; 3D Conv Net (3D-CNN)

gsk

- Focus box on atomic coordinates of amino acid; four atom types (C,O,N,S)
- No "hand-engineered" features, i.e. the network determines relevant features from raw data
 b c



Amino acid substitutions: 3D atomic coordinates; Conv Net

- Alternative method:
 - Convert local environment of 3D point into numeric vector
 - Exact structure not preserved







Amino acid substitutions: 3D atomic coordinates; Conv Net

Learn to predict effect of mutations on protein structure (two alternative approaches)



Amino acid substitutions: 3D atomic coordinates; Conv Net

3D Conv nets are superior for predicting amino acid changes

Confusion matrices for predictions of the 20 amino acid microenvironments.

Heatmap: probability of examples of true label i being predicted as label j.







0.7

0.6

0.1

0.4

ló 3

0.1

0.9

0.8

0.7

0.6

0.2

Amino acid substitutions: 3D atomic coordinates; Conv Net



ARG

PRO

ASP

ASP

ALA

ASP

ASP

GLU

ALA

GLU

ASP

VAL

TYR

ILE

VAL

VAL

PHE

LEU

ILE

LEU

ALA

PRO

ASP

GLU

ALA

VAL

SER

SER

GLU

GLU

GLU

ASP

LEU

GLU

VAL

VAL

ILE

ALA

LEU

ILE

GLU

1L23

1L24

129L

130L

131L

1CV6

1L61

1L19

1L20

1L59

1L62

1L57

1P46

1P6Y

1G0Q

1G0L

1QUG

1TLA

234L

1L17

233L

ALA

PRO

THR

SER

SER

MET

ASN

ASP

ASP

ASN

ASP

ASP

ILE

TYR

ILE

VAL

VAL

PHE

LEU

VAL

LEU

3DCNN agrees on which T4 lyzozyme mutants are **destabilizing or neutral** Comparison: predicted vs actual amino acid at given position for wildtype and mutant

Class	Variant	Wild type	Wild	Wild	Mutant PDB ID	Mutant	Mutant
			ITUB	Prodicted	PUBID	1100	Fredicied
Destabilizing	M102K	2LZM	MET	MET	1L54	LYS	MET
Destabilizing	L99G	2LZM	LEU	LEU	1QUD	GLY	LEU
Destabilizing	V149S	2LZM	VAL	VAL	1G06	SER	ILE
Destabilizing	V149C	2LZM	VAL	VAL	1G07	CYS	VAL
Destabilizing	V149G	2LZM	VAL	VAL	1G0P	GLY	VAL
Destabilizing	T157l	2LZM	THR	THR	1L10	ILE	THR
Destabilizing	G156D	2LZM	GLY	GLY	1L16	ASP	GLY
Destabilizing	R96H	2LZM	ARG	ARG	1L34	HIS	ARG
Destabilizing	D92N	2LZM	ASP	ASP	1L55	ASN	ASP
Destabilizing	I3P	2LZM	ILE	ILE	1L97	PRO	ILE
Destabilizing	V87M	2LZM	VAL	VAL	1CU3	MET	PHE
Destabilizing	R96N	2LZM	ARG	ARG	3CDT	ASN	ARG
Destabilizing	R96D	2LZM	ARG	ARG	3C8Q	ASP	ARG
Destabilizing	R96W	2LZM	ARG	ARG	3FI5	TRP	PHE
Destabilizing	R96Y	2LZM	ARG	ARG	3C80	TYR	ARG
Destabilizing	M102L	2LZM	MET	MET	1L77	LEU	ILE
Destabilizing	M106K	2LZM	MET	LEU	231L	LYS	MET
Destabilizing	M120K	21.7M	MET	GLU	232	LYS	TYB
Destabilizion	137	21 ZM	ILE	IF	11 18	TVB	TVB
Sestabilizing	101	ELE!	Service and	a state of the second second	1210	1111	



Representing small drug-like molecules for machine learning

Conventional representations



Molecular structure graph



- SMILES string
 - CC(=O)Oc1cccc1C(=O)O
- Bit vector fingerprint
 - Different methods (MACCS, Morgan,...) → CDK toolkit, Python package RDKIT
 - Fixed length
 - 1 or 0
 - presence and absence of molecular features
 - can be used directly as input for ML

Chemception: Learning chemistry from 2D drawings; Tox prediction



Table 3: Summary of Results for Tox21 trained on Chemception T1 network.

	Tr	UC	Valid	ation	AUC	Test AUC				
nr-ahr	0.825	+/-	0.018	0.779	+/-	0.015	0.800	+/-	0.020	Y
nr-ar	0.843	+/-	0.010	0.797	+/-	0.049	0.757	+/-	0.029	Y
nr-ar-lbd	0.887	+/-	0.034	0.834	+/-	0.046	0.886	+/-	0.014	Y
nr-aromatase	0.801	+/-	0.010	0.759	+/-	0.027	0.799	+/-	0.016	Y
nr-er	0.747	+/-	0.020	0.710	+/-	0.023	0.694	+/-	0.013	Y
nr-er-lbd	0.824	+/-	0.029	0.765	+/-	0.036	0.762	+/-	0.009	Y
nr-ppar-gamma	0.791	+/-	0.038	0.742	+/-	0.025	0.819	+/-	0.015	Y
sr-are	0.724	+/-	0.009	0.702	+/-	0.025	0.654	+/-	0.009	N
sr-atad55	0.841	+/-	0.022	0.759	+/-	0.048	0.776	+/-	0.011	Y
sr-hse	0.776	+/-	0.032	0.732	+/-	0.013	0.717	+/-	0.018	Ν
sr-mmp	0.791	+/-	0.020	0.759	+/-	0.016	0.755	+/-	0.010	Y
sr-p53	0.844	+/-	0.034	0.782	+/-	0.036	0.776	+/-	0.011	Y
Tox21	0.808		0.044	0.760		0.035	0.766		0.058	

Goh, G. B., Siegel, C., Vishnu, A., Hodas, N. O. and Baker, N. (2017) 'Chemception: A Deep Neural Network with Minimal Chemistry Knowledge Matches the Performance of Expertdeveloped QSAR/QSPR Models', pp. 1–38. Available at: http://arxiv.org/abs/1706.06689.

Generating novel compounds: recurrent neural networks

gsk

Example input molecules with SMILES

- Train RNN (recurrent neural network) model on SMILES strings from ChEMBL (1.4 M molecules)
- 72 M SMILES characters from vocabulary of 51 different characters (one-hot encoded)
- Apply filters to check that output is valid SMILES and drug-like properties (filters)



Segler, M. H. S., Kogej, T., Tyrchan, C. and Waller, M. P. (2017) 'Generating Focussed Molecule Libraries for Drug Discovery with Recurrent Neural Networks', pp. 1–17. Available at: http://arxiv.org/abs/1701.01329.

gsk

Generating novel compounds: recurrent neural networks

- Properties of novel molecules (848000):
 - Phys-chem descriptors very similar to ChEMBL
 - 75% suitable for high-throughput screen in Pharma
 - But new scaffolds (core molecular structure) proposed

Valid SMILES emerge over training time

Batch	Generated Example	valid
0	Oc.BK5i%ur+7oAFc7L3T=F8B5e=n)CS6RCTAR((OVCp1CApb)	no
1000	OF=CCC2OCCCC)C2)C1CNC2CCCCCCCCCCCCCCCCCCCCCC	no
2000	O=C(N)C(=O)N(cloccclOC)c2ccccc2OC	yes
3000	0=C1C=2N(c3cc(ccc30C2CCC1)CCCc4cn(c5c(C1)cccc54)C)C	yes



Segler, M. H. S., Kogej, T., Tyrchan, C. and Waller, M. P. (2017) 'Generating Focussed Molecule Libraries for Drug Discovery with Recurrent Neural Networks', pp. 1–17. Available at: http://arxiv.org/abs/1701.01329.

Generating novel compounds: recurrent neural networks

- Transfer learning:
 - Fine-tuning can be applied to create target-specific predictors
 - Take already trained model (1.4M from ChEMBL) and re-train on known actives for target protein of interest

Active molecules for specific target re-discovered after few additional training epochs with pre-trained model



Table 2 Reproducting known actives in the *Plasmodium* test set.EOR: Enrichment over random.

#	pIC_{50}	Train.	Test	Gen. mols.	Reprod.	EOR
1	> 8	1239	1240	128,256	28%	66.9
2	> 8	100	1240	93,721	7%	19.0
3	> 9	100	1022	91,034	11%	35.7

Table 3 Reproducting known

actives in the Staphylococcus test set. EOR: Enrichment over random.

Entry	<i>р</i> МІС	Train.	Test	Gen. mols.	Reprod.	EOR
1	> 3	1000	6051	51,052	14%	155.9
2	> 3	50	7001	70,891	2.5%	21.6
3^a	> 3	50	7001	85,755	1.8%	6.3
4^b	> 3	50	7001	285	0%	
5^c	> 3	0	7001	60,988	6%	59.6

^{*a*}Fine-tuning learning rate = 10^{-4} . ^{*b*}No Pretraining. ^{*c*}8 Generate-Test cycles.

Segler, M. H. S., Kogej, T., Tyrchan, C. and Waller, M. P. (2017) 'Generating Focussed Molecule Libraries for Drug Discovery with Recurrent Neural Networks', pp. 1–17. Available at: http://arxiv.org/abs/1701.01329.



gsk

Generating novel compounds: Adversarial autoencoders

Train on 6252 compounds profiled in MCF-7 cell lines (breast cancer)



Kadurin, A., Aliper, A., Kazennov, A., Mamoshina, P., Vanhaelen, Q., Khrabrov, K. and Zhavoronkov, A. (2016) 'The cornucopia of meaningful leads: Applying deep adversarial autoencoders for new molecule development in oncology', *Oncotarget*, 5(0). doi: 10.18632/oncotarget.14073.

Generating novel compounds: Adversarial autoencoders





Kadurin, A., Nikolenko, S., Khrabrov, K., Aliper, A. and Zhavoronkov, A. (2017) 'DruGAN: An Advanced Generative Adversarial Autoencoder Model for de Novo Generation of New Molecules with Desired Molecular Properties in Silico', *Molecular Pharmaceutics*, 14(9), pp. 3098–3104.



How does the small molecule bind to the target protein?

DL-based scoring function of binding

gsk

- Again 3D convolution network
- focus on binding site
- Learn to score the binding interaction
- Compare against physics-based scoring function (AutoDock Vina)

Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J. and Koes, D. R. (2017) 'Protein-Ligand Scoring with Convolutional Neural Networks', *Journal of Chemical Information and Modeling*, 57(4), pp. 942–957. doi: 10.1021/acs.jcim.6b00740.

DL-based scoring function of binding



Input: 3D grid 24x24x24 Ligand | Receptor Type 34 atom type channels AliphaticCarbonXSHydrophobe Y Υ AliphaticCarbonXSNonHydrophobe Υ Υ AromaticCarbonXSHydrophobe Y Υ Density distribution around AromaticCarbonXSNonHydrophobe Y Υ atom center Bromine Υ Ν Calcium Ν Υ Chlorine Y Ν Fluorine Υ Ν Iodine Υ Ν Iron Ν Υ Magnesium Ν Υ Nitrogen Υ Υ Y Υ NitrogenXSAcceptor NitrogenXSDonor Υ Υ NitrogenXSDonorAcceptor Υ Υ Y Oxygen Ν OxygenXSAcceptor Υ Υ OxygenXSDonorAcceptor Y Υ Phosphorus Υ Υ Y Sulfur Υ SulfurAcceptor Ν Υ Ν Υ Zinc

Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J. and Koes, D. R. (2017) 'Protein–Ligand Scoring with Convolutional Neural Networks', Journal of Chemical Information and Modeling, 57(4), pp. 942–957. doi: 10.1021/acs.jcim.6b00740.

DL-based scoring function of binding



Atom importance: which parts of the molecule and protein contribute most to binding score (strategy: modify input to understand contribution to output)



Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J. and Koes, D. R. (2017) 'Protein–Ligand Scoring with Convolutional Neural Networks', *Journal of Chemical Information and Modeling*, 57(4), pp. 942–957. doi: 10.1021/acs.jcim.6b00740.

DL-based scoring function of binding

across targets

1.0

0.8

True Positive Rate

0.2

0.2

predictions

0.4

0.6

False Positive Rate

0.8



CNN wins at predicting "good" vs "bad" poses 3MYG Vina CNN **Example: CNN loses** CNN (AUC=0.815) Vina (AUC=0.645) 1.0 CNN loses to Vina at same-target 3PE2 Vina CNN

Example: CNN wins

Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J. and Koes, D. R. (2017) 'Protein-Ligand Scoring with Convolutional Neural Networks', Journal of Chemical Information and Modeling, 57(4), pp. 942-957. doi: 10.1021/acs.jcim.6b00740.

Calculating binding free energy



- Characterize molecular neighborhood of each atom (distances to nearby atoms, atom types)
- Learn/predict energies of bound and unbound states → free energy (strength of drug binding)



Gomes, J., Ramsundar, B., Feinberg, E. N. and Pande, V. S. (2017) 'Atomic Convolutional Networks for Predicting Protein-Ligand Binding Affinity', pp. 1–17. Available at: http://arxiv.org/abs/1703.10603.

Calculating binding free energy



- Characterize molecular neighborhood of each atom (distances to nearby atoms, atom types)
- Learn/predict energies of bound and unbound states → free energy (strength of drug binding)





Calculating binding free energy



- Characterize molecular neighborhood of each atom (distances to nearby atoms, atom types)
- Learn/predict energies of bound and unbound states → free energy (strength of drug binding)

Compare against known free energies

	ACNN		GRID-RF		GRID-NN		GCNN		ECFP-RF		ECFP-NN	
Split	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
Random	.912	.448	.969	.336	.963	.058	.676	.265	.920	.212	.942	.227
Stratified	.939	.116	.969	.132	.963	.165	.735	.064	.924	.071	.942	.077
Scaffold	.911	.043	.965	.109	.953	.067	.797	.254	.920	.218	.940	.206
Temporal	.923	.251	.972	.287	.957	.245	.744	.095	.925	.206	.952	.071

Table 1. Performance (Pearson R^2) on PDBBind core train/test sets.

	ACNN		GRID-RF		GRID-NN		GCNN		ECFP-RF		ECFP-NN	
Split	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
Random	0.325	0.774	0.385	0.741	0.230	0.877	0.656	1.112	0.399	1.112	0.234	1.138
Stratified	0.282	0.997	0.339	0.990	0.205	0.813	0.556	0.995	0.410	0.901	0.223	1.115
Scaffold	0.410	0.993	0.338	1.397	0.211	1.630	0.516	0.883	0.438	1.003	0.221	0.909
Temporal	0.363	0.974	0.368	0.860	0.237	0.809	0.588	1.062	0.413	0.974	0.341	1.265

Uses DeepChem...

 Table 2. Performance (MUE [kcal/mol]) on PDBBind core train/test sets.

Gomes, J., Ramsundar, B., Feinberg, E. N. and Pande, V. S. (2017) 'Atomic Convolutional Networks for Predicting Protein-Ligand Binding Affinity', pp. 1–17. Available at: http://arxiv.org/abs/1703.10603.

DeepChem

Deep Learning tools for drug discovery

gsk

- Vijay Pande lab (Stanford)
- Implementation of many useful deep learning techniques for small molecules / drug binding (e.g. Graph convolution, dataset stratification based on molecular scaffold, ...)



https://www.deepchem.io/

deepchem 1.3.

Tutorials

The following tutorials show off various aspects or capabilities of DeepChem. They can be run interactively in Jupyter (IPython) notebook. Download the notebook files and open them in Jupyter:

\$ jupyter notebook

DeepChem Tutorials

- Predicting Ki of Ligands to a Protein
- · Conditional Generative Adversarial Network
- · Creating a high fidelity DeepChem dataset from experimental data
- Graph Convolutions For Tox21
- · MNIST with DeepChem and TensorGraph
- MNIST GAN
- Multitask Networks On MUV
- Pong in DeepChem with A3C
- Basic Protein-Ligand Affinity Models
- Tutorial: Use machine learning to model protein-ligand affinity.
- The protein-ligand complex view.
- Quantum Machinery with gdb1k
- SeqToSeq Fingerprint
- Modeling Solubility
- TensorGraph Mechanics

Contributing tutorials

Do you have a neat example of using DeepChem? Format your code into an IPython notebook and submit a pull request!

Source

© Copyright 2016, Stanford University and the Authors. Created using Sphinx 1.3.5.

41



Complex effects of drugs inside an organism

specific toxicological effects.



FIGURE 6 | DeepTox pipeline for toxicity prediction.

Mayr, A., Klambauer, G., Unterthiner, T. and Hochreiter, S. (2016) 'DeepTox: Toxicity Prediction using Deep Learning', *Frontiers in Environmental Science*, 3(February). doi: 10.3389/fenvs.2015.00080.

CGL

MML

NCI

VIF

DeepTox: Toxicity Prediction using Deep Learning



Mayr, A., Klambauer, G., Unterthiner, T. and Hochreiter, S. (2016) 'DeepTox: Toxicity Prediction using Deep Learning', Frontiers in Environmental Science, 3(February). doi: 10.3389/fenvs.2015.00080.



Predicting tox based on molecular features, trained on clinical trial outcomes Highlights

- Computational approach predicts the likelihood of clinical trial toxicity
- Identification of molecule and target properties associated with clinical toxicity
- Development of a tool to facilitate interaction and interpretation of the model





Gayvert, K. M., Madhukar, N. S. and Elemento, O. (2016) 'A Data-Driven Approach to Predicting Successes and Failures of Clinical Trials', *Cell Chemical Biology*. Elsevier Ltd, pp. 1–8. doi: 10.1016/j.chembiol.2016.07.023.





Predicting tox based on molecular features, trained on clinical trial outcomes



Gayvert, K. M., Madhukar, N. S. and Elemento, O. (2016) 'A Data-Driven Approach to Predicting Successes and Failures of Clinical Trials', *Cell Chemical Biology*. Elsevier Ltd, pp. 1–8. doi: 10.1016/j.chembiol.2016.07.023.



Predicting tox based on molecular features, trained on clinical trial outcomes





Gayvert, K. M., Madhukar, N. S. and Elemento, O. (2016) 'A Data-Driven Approach to Predicting Successes and Failures of Clinical Trials', *Cell Chemical Biology*. Elsevier Ltd, pp. 1–8. doi: 10.1016/j.chembiol.2016.07.023.

Review articles

Deep Learning in Bio/Chem/Med/Pharma



- Mamoshina, P., Vieira, A., Putin, E. and Zhavoronkov, A. (2016) 'Applications of Deep Learning in Biomedicine', *Mol Pharm*, 13(5), pp. 1445–1454. doi: 10.1021/acs.molpharmaceut.5b00982.
- Min, S., Lee, B. and Yoon, S. (2016) 'Deep Learning in Bioinformatics'. doi: 10.1093/bib/bbw068.
- Angermueller, C., Pärnamaa, T., Parts, L. and Stegle, O. (2016) 'Deep learning for computational biology', *Molecular Systems Biology*, 12(7), p. 878. doi: 10.15252/msb.20156651.
- Gawehn, E., Hiss, J. A. and Schneider, G. (2016) 'Deep Learning in Drug Discovery', *Molecular Informatics*, 35(1), pp. 3–14. doi: 10.1002/minf.201501008.
- Baskin, I. I., Winkler, D. and Tetko, I. V. (2016) 'A renaissance of neural networks in drug discovery.', *Expert opinion on drug discovery*. Taylor & Francis, 441(June), p. 17460441.2016.1201262. doi: 10.1080/17460441.2016.1201262.
- Goh, G. B., Hodas, N. O. and Vishnu, A. (2017) 'Deep learning for computational chemistry', *Journal of Computational Chemistry*, 38(16), pp. 1291–1307. doi: 10.1002/jcc.24764.
- Chen, Y., Li, Y., Narayan, R., Subramanian, A. and Xie, X. (2016) 'Gene expression inference with deep learning', *Bioinformatics*, 32(12), pp. 1832–1839. doi: 10.1093/bioinformatics/btw074.
- Hodos, R. A., Kidd, B. A., Shameer, K., Readhead, B. P. and Dudley, J. T. (2016) 'In silico methods for drug repurposing and pharmacology', *Wiley Interdisciplinary Reviews: Systems Biology and Medicine*, 8(3), pp. 186–210. doi: 10.1002/wsbm.1337.
- Shen, D., Wu, G. and Suk, H.-I. (2017) 'Deep Learning in Medical Image Analysis', *Annual Review of Biomedical Engineering*. Annual Reviews , 19(1), pp. 221–248. doi: 10.1146/annurev-bioeng-071516-044442.



Thanks





heidelberg.ai

Dec. 12th, 6pm: [Paper Discussion] Dynamic Routing between Capsules