

Talk 1 - Y. Bengio - Deep Learning I

28/10/17

1. Backprop

2. DNNs

3. RNNs

} Broad intro
together afterwards

- attention mechanisms → liberate NNs from fixed-size input
 - ↳ "flow is a dog"
 - ↳ between less attention at certain spot of image while spelling.



- powerful partitions → can benefit curse of dimensionality
 - ↳ mimic how brain works with semantics
 - ↳ compositionality assumption: Exponential gain in representational power
 - ① ⇒ distr. representations
 - ② ⇒ deep architecture (layers of features) → sequential
 - ↳ Assumption: Composites are able to grasp underlying phenomena in the world!

① logical partitions vs. distr. partitions

② world is compositional

- attention ⇒ gating units → control info flow!

↳ E.g. neural translation

↳ System I: Gtex → fact → intuition } controller
System II: Hippocampus → store & recall } to read and write

- SGD → extremely simple method

↳ But: extremely powerful technique!

- Problems
 - Current models cheat by picking on specific regularities → superficial depth
 - Not able to capture long-term dependencies → long-scale dependency
 - Heavily relies on smooth differential functions
- Autoencoders
 - Mapping of low & high-levels of abstraction
 - distortion of space → deform complex dist. into space where it is simple to do computations → manifolds disentanglement
 - ↳ 0-shot learning: learn without data → prob. mapping between representation spaces and embedded parts => classification based on similarity in different embedded spaces
- Challenges
 - Unsupervised learning
 - Meaningful features \Leftrightarrow causal features of the world
 - Connection to RL \Rightarrow reason about rare events
 - Causal framework to causal reasoning
 - Hardware: specialized chips/GPUs for MLs
 - Representation learning: "Good" representation \Leftrightarrow indep. controllable of the world
 - ↳ discovery by acting in the world

Talk 2 - S. Meyn - RL

28/08/17

- RL \subset Stoch. Approx.

$$\rightarrow \text{Stoch. Approx.: } \hat{f}(\theta^*) := \left[E[f(\theta, w)] \right]_{\theta=\theta^*} = 0$$

exp. \Rightarrow complete
refinement

potentially complex / unknown

\rightarrow Monte Carlo Average Estimation: $\xrightarrow{\text{ergodicity}}$

$$\theta(u) = \frac{1}{n} \sum_{i=1}^n c(X(i)) \xrightarrow{\text{recursion available}}$$

\rightarrow MDPs: $\rho \{ X(t+1) \in A | X(t), U(t) = u, \text{prev hist.} \}$

decision u

↳ $h^*(x) = \min_u \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[c(X(t)), U(t)] | X(0) = x]$

discount cost

↳ satisfies Bellman \rightarrow Galekin relaxation available

↳ Q-learning: Fixed point eq. for a pt. [Watkins]

\Rightarrow with outside of expectation \rightarrow stochastic Opt.

↳ Ergodicity does not have to hold in reality: e.g. stocks!

↳ In practice: slow and unstable

\rightarrow recursion: $\theta(t+1) = \theta(t) + \alpha \underbrace{f(\theta(t), X(t))}_{\text{gain-critical!}} \xrightarrow{\text{approx.}} \theta^*$

Fuler
approx.
policy
② (nice)

ODE: $\frac{d}{dt} \varphi(t) = \hat{f}(\varphi(t)) \rightarrow \theta^*$

\rightarrow forhar's Monograph \Rightarrow best reference

④ valuable relationship!

- Performance Eval.
 - Finite-d_n bound \Rightarrow concentration!
 - Asymptotic variance \rightarrow not really applicable for anything
 - \hookrightarrow MC approx. \Rightarrow No concentration result; beautiful asymptotics
 - \hookrightarrow Lyapunov Eq. \Rightarrow not used even though it is stronger than bound! \rightarrow problem with bad convergence
- How to model the gain / learning rate?
 - Least squares TD
 - inversely proportional to optimal variance
 - $\Sigma^* = \hat{\Sigma}^{-1} \sum_{\Delta} \hat{\Sigma}^{-1}$
 - \rightarrow Polyak-Ruppert
 - \rightarrow Stochastic Newton-Raphson
- $\exists \Delta \in \Theta$ learning
 - $\Delta \in \Theta \xrightarrow{?} \Delta^*$ \Rightarrow ODE: Newton-Raphson with nonlinearities!
 - $c(\Delta) \rightarrow c^* \Leftrightarrow \Delta \rightarrow \Delta^* \rightarrow 1-\text{to}-1 \text{ mapping!}$
- RL Open Qs // Problems
 - solved by dimensionality reduction
 - finite concentration results \rightarrow have to be used to design algorithm design!
 - Context-specific algo designs
- Trade-off between c in $\Theta(c)$ and concentration under the covariance metric.

Talle 3 - 4. Begie - Deep Learning II

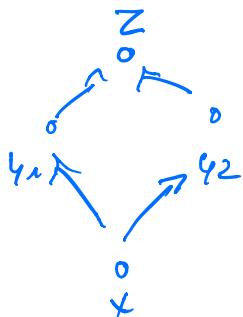
28/08/17

① Backprop and graph flow

- Objective function is not underlying truth (\Rightarrow data-dependent) \rightarrow don't have to find exact minimum! \rightarrow approx. opt.
- Gradient \Rightarrow best direction to go to \rightarrow look at proof!
- Forward pass: $O(n)$ \Rightarrow where $n = \#$ neurons
- Backprop: $O(n)$ \rightarrow takes $O(n^2)$
- Alternative: Finite differences

$$\frac{\partial L(\theta_i, \theta_{-i})}{\partial \theta_i} = \frac{L(\theta_i + \epsilon, \theta_{-i}) - L(\theta_i, \theta_{-i})}{\epsilon}$$

- backward accumulation: Only gradient computation method
- \rightarrow chain rule: Consider intermediate values
to learn how to change parameters by learning how to change intermediate values!



$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x} \\ \frac{\partial z}{\partial x} &= \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}\end{aligned}$$

- Loss fct.: negative log LL $- \log p(y|x=x) \rightarrow$ cross-entropy!
- $\frac{\partial z}{\partial x}$ interesting if we want to measure sensitivity/robustness!

- Very general framework! \rightarrow just have to define forward prop \Rightarrow implementations define the these operations needed!
- Automatic Differentiation \rightarrow gradient comp. can be automatically inferred from symbolic expression of f prop.
- Neural justification
 - \rightarrow earlier: people were sceptical since forward and backwords computation is very different
 - \hookrightarrow new findings suggest that brain is able to use neurons in both stages
 - \rightarrow new neural findings! \Rightarrow makes it more reasonable!
- Property: Change in graph \rightarrow stochastic generalization
- Theano, Tensorflow: Work on static graphs \rightarrow can adapt to dynamic sequences \Rightarrow but can break down
 \hookrightarrow pytorch: does construct graphs on the fly!
 - Higher-level abstractions disentangle the factors of variation!
 - Good representation?
 - Smoothness prior / assumption \rightarrow often turns not accurate for high-dimensional spaces! \Rightarrow make function smooth by further construction

Wifi Ku: 356 E 8

Talk 4 - Y. Bengio - Deep Learning III

29/08/17

- Hidden units \Rightarrow Meaning is not directly clear ② DNNs
 - \rightarrow learning "discovers" significant features
 - \rightarrow No hand engineering / pre-defined meaning for hidden units
- Universal approx. Th.: single hidden layer
 - \rightarrow Why deep? \Rightarrow Functions can be approx. better / easier (due to compositionality) with less parameters!
 - \rightarrow Problems:
 - * reasonably deep net vs. huge shallow net
 - * not early optimization guarantee
 - * not good generalization guarantee (finite training sets)
 - \Rightarrow Only tells us existence \rightarrow does not guarantee that we find weight constellation
- breakthroughs \Rightarrow ReLU \rightarrow makes training possible
 - \rightarrow composition of nonlinearities \Rightarrow no vanishing gradient
 - \rightarrow sigmoid \Rightarrow generic unit
 - \rightarrow tanh \Rightarrow used a lot in RNNs?
 - \rightarrow softmax \Rightarrow sometimes used in attention \rightarrow not an elementwise function \rightarrow involves interaction between neurons \Rightarrow competition between nearby neurons
 - \rightarrow ReLU can't be used in RNNs \Rightarrow not upper bound
to values can blow up
- \hookrightarrow "symmetry breaking" \rightarrow Never-off \Rightarrow no weight decay
 - \rightarrow Never-on \Rightarrow strong decay \rightarrow specialized networks
- \hookrightarrow Gradient seeds credit information \Rightarrow no diffusion of information! \rightarrow no solid explanation

- Expert knowledge \rightarrow does not get lost \Leftrightarrow substitution with new architecture knowledge
- Bayesian NNs \rightarrow very good on small datasets where overfitting is a problem \rightarrow sample from posterior gives uncertainty
 \rightarrow not a good application yet!
- RelU problem \rightarrow some neurons might die \rightarrow negative weights
↳ such \rightarrow never get act
↳ Furthermore: Only a few neurons are exhibiting RelUs
- Noise injection in sigmoid activation fn. to help not getting stuck with the gradient!
- Loss function: Negative log-Ltt \Leftrightarrow cross-entropy

$$-\log f(x)_y = -\log p_{\theta}^{f(x)}(y|x)$$

↳ log: simplifies numerical stability ad with
- $y \in \mathbb{R}^d$ cond. given X
- indep.: $-\log P(Y|X) = -\sum_i \log P(Y_i|X)$
add up the losses
↳ missing values in $y \rightarrow$ free loss!
- dep.: need to capture cond. joint dist.

$$P(Y_1, \dots, Y_d | X)$$

↳ use factorizable to sequentially predict with max-Ltt framework

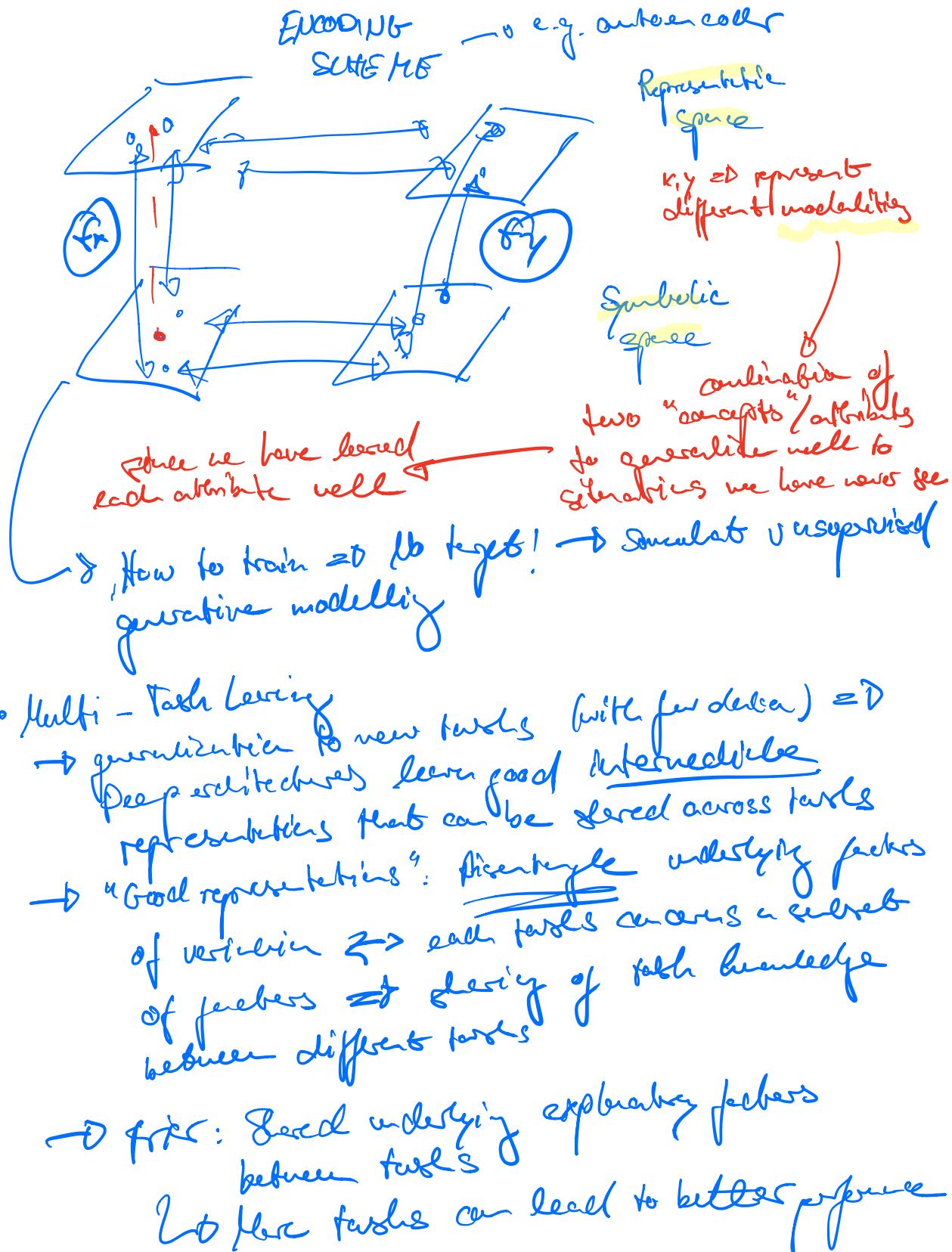
-
- (cont.) afterwards use sigmoid fn. for best prediction that is hard (high var)
 - right now not applicable
↳ small datasets!

Final Thought

- Bayesian VI for NNs
↳ FTFP \rightarrow model computes
- new perspective!
- new comp. tools // question the modelling perspective

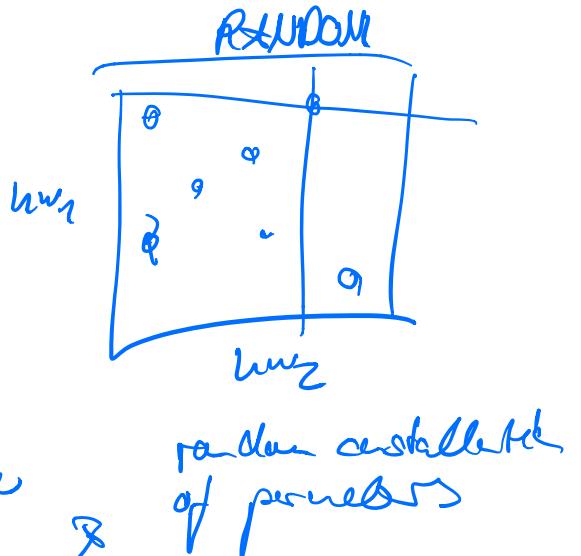
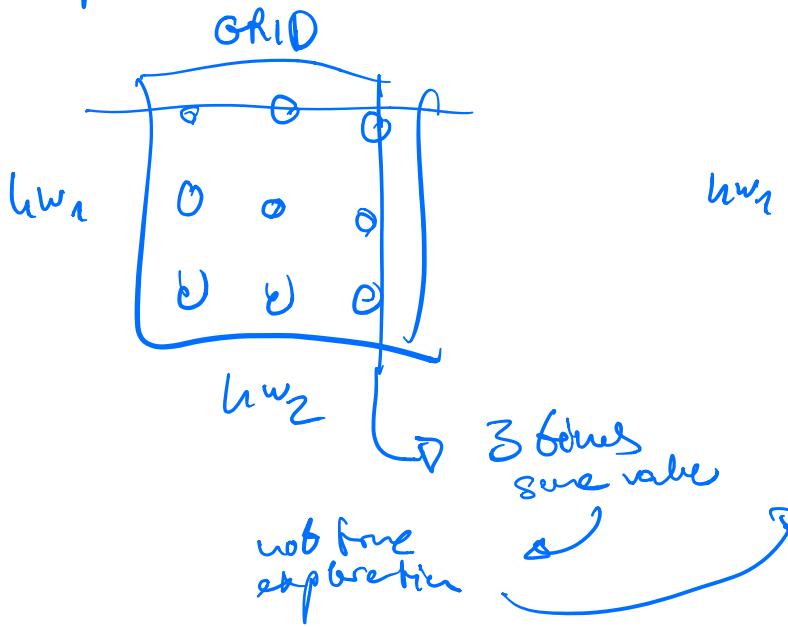
benign

- ELBO is biased
Train - Val - Test losses
- Use train. framework to construct first linear representations
↳ weights have small variance



- Hyperparameters: Design decisions \oplus opt. stability
 - \rightarrow Train \leftrightarrow Validability \leftrightarrow Test
 - \downarrow
 - Form of optimizable!
 - hyperparam \neq generalizability
train error \neq test error
 - \hookrightarrow shuffle and split dataset \Rightarrow no assumption which
 - \rightarrow most important parameter: global learning rate
 - \hookrightarrow in some sense equivalent to choosing the learning algorithm!
 - \rightarrow minibatch: parallel GPU computation
 - \rightarrow AdaDelta / Adagrad \rightarrow allows to combat adversarial examples! (\hookrightarrow no risk of too much noise)
 - \rightarrow cross-validation: makes sense if dataset is small!
- Sequential Evaluation: TS dataset
 - \nearrow purpose
 - TV & \rightarrow simulate what would have happened
 - TTVB \rightarrow if we had dataset only up to t
 - TTTVB \rightarrow \hookrightarrow repeat with reversed f and average
 - double CV \hookrightarrow fast performance over all labels
 - Grid search \Rightarrow inefficient with more than 2 hyperparameters
 - Random search \Rightarrow simple \oplus efficient
 - Bayesian optimisation \rightarrow still need to set interval of values for parameters

- Random search



- Non penalization

→ groups penalization:

$$\mathcal{E}(\theta) = \sum_i Q_i \gamma_i(\theta) + \lambda \left[\sum_j w_j^2 \right]^{\frac{1}{2}}$$

↑ eliminate whole group of weights

↑ per row selection

- weight initialization: Glorot initialization

→ biases: 0

→ weights: factor \approx propagation of info through gradients

$\hookrightarrow \neq 0$ for tank \Rightarrow or slow that all gradients 0

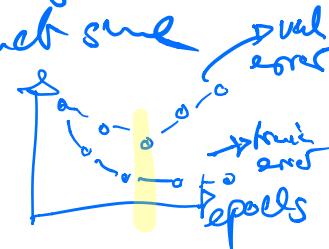
$\hookrightarrow \neq 0$ for tank \Rightarrow middle units will act slow

$\hookrightarrow \neq$ one value \Rightarrow middle units will act slow

\Rightarrow simple $w_{i,j}^{(l)}$!!

- Early stopping: Under of training iterations

→ store training and validation error and choose sweet spot



Talk 5 - P. Ravinder - GMs I

28/08/17

① Representation

- GMs and representable problems are complementary!
- $\exists \rightarrow$ Effective reasoning in complex environments:
 - I. Declarative Knowledge: Structured knowledge base - e.g. prior
 - II. Reasoning: derive optimal output to solve inference problem

↳ problem with logical reasoning: uncertainly
knowledge base certainty

↳ GM provide framework to perform probabilistic reasoning in a comp. tractable way!

- GMs models declarative knowledge as a multivariate dist. over random variables \rightarrow highly complicated function!

- (V, E) : Each node $i \in V$ has a set of parents $\overline{\pi}_i$
↳ DPG: $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | \overline{x}_{\pi_i})$ joint distribution

\rightarrow storage cost of general joint distr. where each var. can take k values: Γ^n

\rightarrow storage cost of a DPG: for var. i , $\Gamma^{|\pi_i|}$

↳ as long as nodes don't have many parents (fan-in) we are able to compactly store the joint dist. \rightarrow go from global to local

- Undirected graph \Rightarrow loose symmetry with directed edges
 \Rightarrow Hidden Markov Model \rightarrow FEEDBACK
- \rightarrow notion of locality needed for compact representation
- \Rightarrow neighbors \rightarrow consistency of conditional dist.
- \hookrightarrow maximal clique \Rightarrow set of fully connected nodes
- \hookrightarrow potential function $\phi_c(x_c)$: only depends on values x_c of max. clique variables X_c .

\star non-negative, real-valued, otherwise arbitrary

$$p(x) := \frac{1}{Z} \prod_{c \in C} \phi_{X_c}(x_c)$$

set of max. cliques

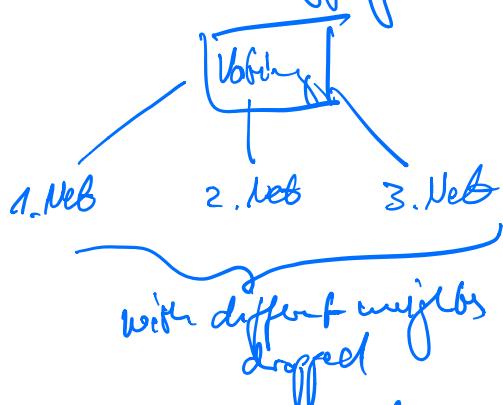
$$Z := \sum_x \prod_{c \in C} \phi_{X_c}(x_c) \rightarrow \begin{array}{l} p \\ \text{normalization constant} \end{array}$$

\hookrightarrow "partition function" \Rightarrow hard to compute!

Talk 6 - Y. Bengio - Deep Learning IV

22/08/17

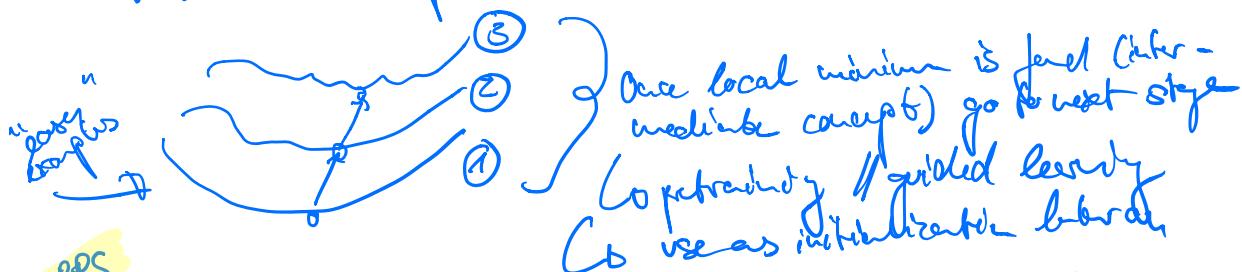
- Proport \Rightarrow Efficient Bagging \rightarrow Indirect ensemble method



only at testing time
to learn to be robust to
noise injection!

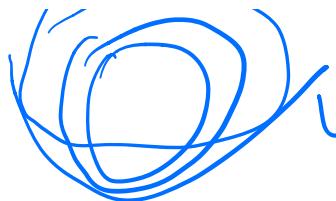
- Underfitting \Rightarrow better optimization
- Overfitting \Rightarrow use data / better regularization

- Conclusions becoming
 \rightarrow learn models sequentially! \Rightarrow family of objective functions



TIPS

- Train on small datasets \Rightarrow know that if we perform too well that there might be a bug! \rightarrow adopt code!
- Watch behaviour of error curves during training!
 \rightarrow is the step right?
- Compare with simple baseline models \rightarrow does it always work?
- Look at examples that lead to overconfident networks
 \rightarrow logit regression
- There is no percentage for train-test \Rightarrow just the right split that makes things smooth!

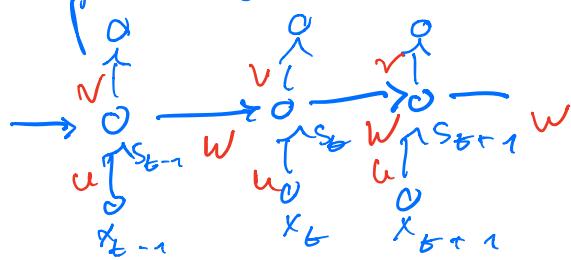


→ King of learning rate as a ball!
smaller l.r. makes ball smaller!

- reduced variance optimization methods → then unclear how to apply!

③ RNNs

- recursive v. state: $s_t = f_\theta(s_{t-1}, x_t)$

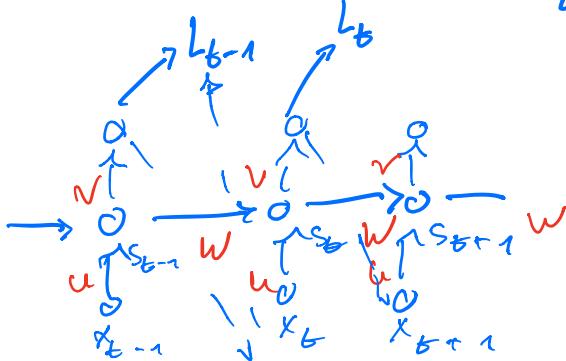


flexible output/input
length
↳ unfold to obtain
backprop structure

- RNN as directed generative model!

$$P(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | x_{t-1}, \dots, x_1)$$

$$L_\theta = -\log P(x_T | x_{T-1}, \dots, x_1)$$



S2V
S2Sfixed
V2S
S2S

- "Teacher forcing" \leftrightarrow MLE

→ RESULT

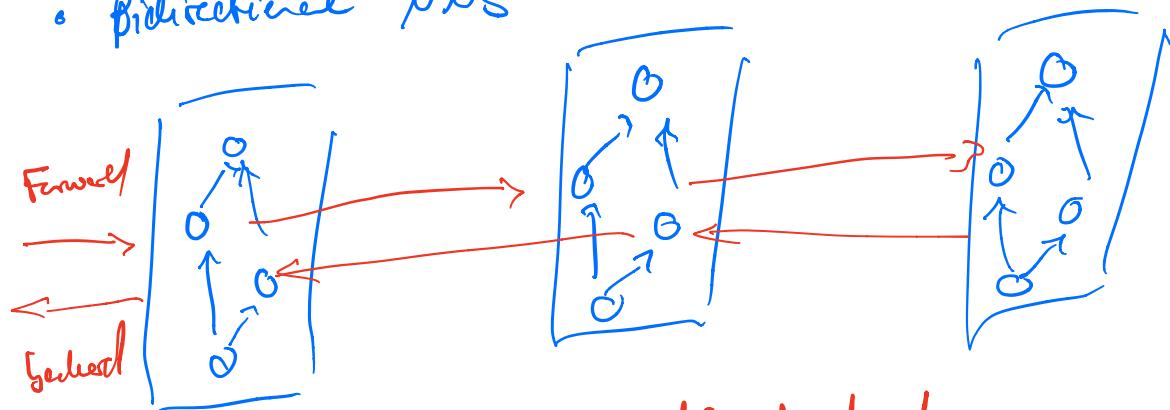
→ training vs. usage: input conditions at training time might be very different compared with test time

↳ solution: inject noise during training

metaphor:
driving or
autonomously

→ "profess" → use GRUs

- Skip connections: bypass layer! ⇒ shorter paths
- Bidirectional GRUs



↳ backprop only requires DTG structure!

→ need to know the full sequence!

- Multiplicative interactions (dot products) seem to work well → febby
- long-term relations ⇒ dynamical systems!
- Variability gradient problem!
- different for MLP and RNNs ⇒ height slicing

$$\frac{\partial L}{\partial s_g} = \frac{\partial L}{\partial s_f} \frac{\partial s_f}{\partial s_{g-1}} \dots \frac{\partial s_{g+1}}{\partial s_g}$$

→ $\lambda(\text{Jacobian}) \geq 1$: gradient explodes → can't trust
 $\lambda(\text{Jacobian}) \leq 1$: gradient vanishes → large gradients!

$$\rightarrow \text{gradient norm clipping} \Rightarrow \text{constant gradient}$$

$$\hat{g} \leftarrow \frac{\text{grad}}{\|\hat{g}\|} \quad \hat{g} \Rightarrow \text{gradient explosion!}$$

- LSTMs \Rightarrow Building unit \rightarrow sigmoid \Rightarrow input, forget, output

\Rightarrow forget mechanism \Rightarrow reading writing in neural net machine!

\hookrightarrow forgetting: vanishing gradients

\hookrightarrow memory: higher-dim. state \Rightarrow avoid need for forget

\hookrightarrow memory: higher-dim. state \Rightarrow avoid need for forget

- different time layers! \Rightarrow pyramidal layer state design!
- delays \Rightarrow skip connections over time!

New Situations \rightarrow search for underlying truth!

Given more abstract phenomena of world

\hookrightarrow Good representability

\rightarrow controllable features \Rightarrow want to have jointly representable

and policies! \Rightarrow latent properties can be controlled

indep. from outer flags in the environment!

\hookrightarrow combine RL with deep learning

\Rightarrow above inference $\xrightarrow{\text{FEP}}$

Tatjana - P. Kaufler - GMs II

29/08/17

- Partition functions \Rightarrow cannot be cond. probabilities \rightarrow do not integrate to 1!

\hookrightarrow also: marginal probs don't work! \Rightarrow "double entry"

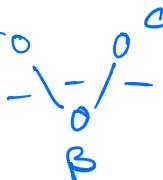
\rightarrow pre-probabilistic interpretation

- Distributions are constrained in terms of conditional independence

$$X_A \perp\!\!\!\perp X_B : p(X_A, X_B) = p(X_A) p(X_B) \quad \text{indep.}$$

$$X_A \perp\!\!\!\perp X_C | X_B : p(X_A, X_C | X_B) = p(X_A | X_B) p(X_C | X_B) \quad \text{cond. indep.}$$

\hookrightarrow general joint is a fully connected undirected graph!
 ≥ 0 no cond. indep. assumption



$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \quad \boxed{\text{joint}}$$

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | \bar{x}_{\neq i}) \quad \boxed{\text{post.}}$$

- graph allows us to formulate conditional independence ass. and to obtain special factorization of the joints.

$$\begin{array}{ccc} x & \xrightarrow{q} & z \\ o \rightarrow o \rightarrow o \end{array}$$

$$p(z|x,y) = p(z|y)$$

$$\Leftrightarrow x \perp\!\!\!\perp z | y$$

$$\begin{array}{ccc} x & \searrow & z \\ & o & \\ & \downarrow & \\ & y & \\ & \downarrow & \\ x \perp\!\!\!\perp z \end{array}$$

$$\begin{array}{ccc} y & \searrow & z \\ & o & \\ & \downarrow & \\ & z & \\ & \downarrow & \\ x \perp\!\!\!\perp z | y \end{array}$$

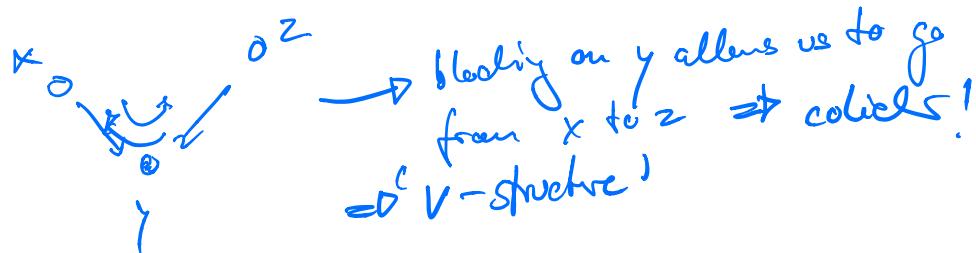
But not $x \perp\!\!\!\perp z | y$!

- d-separation algorithm
 - blocking rules

$$x \rightarrow y \rightarrow z$$

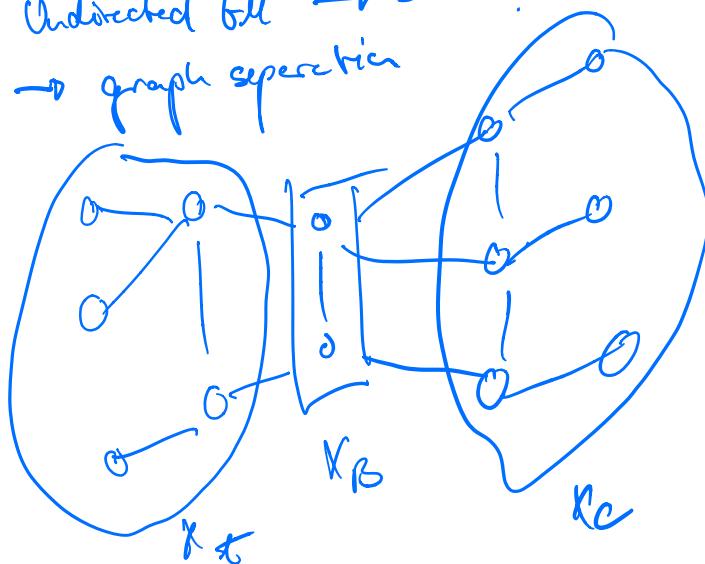
if I condition on y,
I can't go to z!

$$x \rightarrow y \rightarrow z$$



- Undirected fill \rightarrow easier!

\rightarrow graph separation



= advantage:
but factors don't have simple interpretation!

$$\nexists x_* \perp\!\!\!\perp x_c \mid x_b$$

Obs 1: Given G, all dist. expressed as product of conditional
prob. over maximal cliques \Rightarrow PCP-B. THEORY

Obs 2: Given G, all dist. that satisfy all cond. indep. assump's
specified by graph separation \Rightarrow ORTPH THEORY

\Rightarrow Equivalence \rightarrow Hammersley Clifford theorem

- Fitted GM:

Def 1: Given G , $p(x_1, \dots, x_n) \triangleq \prod_{i=1}^n p(x_i | K_{x_i})$

Def 2: Given G , all dist. that satisfy all ad. independence assumptions specified by Bayes Ball alg

\Rightarrow Equivalence

Fall 8 - P. Ravinder - GMs III

30/08/17

- Graph \leftrightarrow Locality \leftrightarrow factor functions
- Joint Distr. Representation

- Edges \rightarrow
 - (X) Correlation \rightarrow Conditional Independence
 - (X) Consistency \rightarrow Multiple testing ref?
- Lack of edge \rightarrow independence of two nodes given rest of nodes
- Factor functions \Rightarrow fcts. over maximal cliques \rightarrow undir.

\Rightarrow cond. prob. distr. of nodes given parents \Rightarrow Dir.

To GM: set of distr. give graph structure
 * non-parametric class of distr. \Rightarrow if discrete always
 * practice: specify parametric class to give factor fit. spec. form

- Fising Model \Rightarrow binary RV \rightarrow spin models
 \rightarrow generalization to discrete GMs or categorical RV

Lager cliques lead to more and more complex factor functions \rightarrow focus on small 2-size cliques! \rightarrow Simplify

→ factor function specification can be linear/non-linear/etc.

- Gaussian GM: Variables are fin.-valued continuous

$$p(x, \theta) := \exp \left\{ - \sum_{(s,t) \in E(G)} \theta_{st} x_s x_t + \log \det(\theta) \right\}$$

↳ joint multivariate Gaussian distribution

↳ θ_{st} : scaling of quadratic \Rightarrow relation to precision!
 \Rightarrow nice interpretation of edge weights

- High-dim problem: combat by combat by compact representation of factors!

- Cont value label/ Slotted continuous data

\Rightarrow Exponential family GMs

$$\rightarrow \text{Node-cond. distr.}: p(x_i | x_{-i}) \xrightarrow{\text{all other nodes}}$$

↳ Fstg: Bernoulli

↳ Categorical: Multinomial

↳ GGM: univariate Gaussian

$$\rightarrow p(x_i | \theta) := \exp \left\{ \sum_{i \in I} \theta_i B_i(x) + C(x) - t(\theta) \right\}$$

univariate
distr.
family

params.

affine
sts

base
mvar

log partition
normaliz.

→ existence of consistent joint distr. if each node-cond. distr. is specified by exp.-family distr.

$$p(X) = \exp \left\{ \sum_s \Theta_s \beta_s(X_s) + \sum_{s \in V} \sum_{t \in N(s)} \Theta_{st} \beta_s(X_s) \beta_t(X_t) \right. \\ \left. + \sum_{s \in V} \sum_{t_1, \dots, t_k} \Theta_{s-t_1, \dots, t_k} \beta_s(X_s) \prod_{j=1}^k \beta_{t_j}(X_{t_j}) \right. \\ \left. + \sum_s g(k_s) - k(\Theta) \right\}$$

↳ factor function is just product of sufficient statistics at eligible nodes scaled by Θ ! \Rightarrow tensor prod.

↳ nice compact factor functions \Rightarrow allows us to model unlabelled data!

↳ $k > d$: Is a problem \Rightarrow need more structure!

- Other constraint for cond. independence:
 \rightarrow reduced parametrization \Rightarrow algebraic structure imposed
 \hookrightarrow constraint to lower order factors!



- Goals:
 - ① Learning the graph structure \rightarrow Model Selection
 - ② Learning the parameters of a parametric model \rightarrow Inference

$$p(X; \Theta, G) = \frac{1}{Z(\Theta)} \exp \left\{ \sum_{(s,t) \in E(G)} \Theta_{st} \beta_{st}(X_s, X_t) \right\}$$

↓
library potential function

[I]a) Model Selection \Rightarrow undir. care. \rightarrow deal with normalizing constant
 \Rightarrow NP-hard!

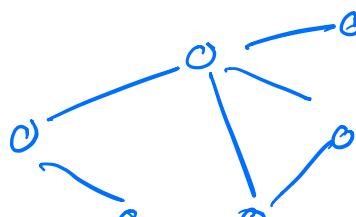
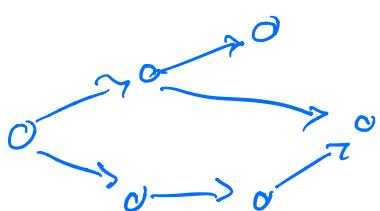
- score-based approaches : search over space of graphs with score per score based on parameter inference
- constraint-based approaches : hypothesis test for cond. independence
 \rightarrow needs a lot of samples if we cond. on a lot of variables \Rightarrow sometimes only constraint on subtrees
 \rightarrow also: multiple testing score correction !
- Estimating neighborhoods : go back to locality !
 \rightarrow high - order stat. model est.
 \rightarrow high - order hyp. testing

\hookrightarrow sparse logistic regression! for each node \rightarrow regularization
 \Rightarrow get neighborhood by counting neighborhoods estimated in a consistent manner (AND, OR)

\hookrightarrow regularization : $d_n \geq \sqrt{\frac{\log p}{n}}$ \Rightarrow ~~overestimates~~ result recovers exact graph structure with very very high prob.

[II]b) Model Selection \Rightarrow directed care

- Moralization \Rightarrow merge all parents \rightarrow fully connected points



learn with graph signal
 "local inside"

- Orient edges using cond. independence tests
 \Rightarrow PC algo \rightarrow Not very scalable !

2] Inference

- Need to approx. \Rightarrow otherwise inference is NP-hard
- We now structure, edge weights and exp.-family distr.

$$p(x_F | x_E) \rightarrow p(x_E, x_F) = \int p(x_E, x_F, R) dx_R$$

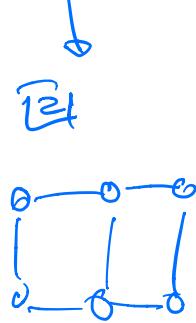
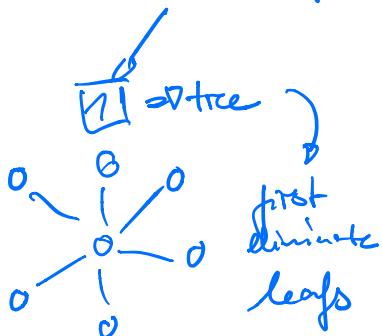
$$= \frac{p(x_E, x_F)}{p(x_E)}$$

$p(x_E) = \int p(x_E, R_F) dR_F$

↑ can take exponential time if R is very large! ↳ high-dim integral

• Exact Algos

- exp. in tree-width of graph
- Variable Elimination \Leftrightarrow leverages locality
- Sum-product for tree-structured graphs
- Sum-product for tree-structured graphs
- ↳ Variable Elimination: variable-wise marginalization
to keep creating new local factors!
- Elimination \Leftrightarrow connector of all its neighbors
- Comp. complexity is exp. in size of largest clique obtained via eliminating variables
- Tree-width: 1 minus size of largest clique size
- Tree-width: 1 minus size of largest clique size
to NP hard to find best possible ordering of eliminate order!



↳ Sum-product:

$$p(x) \propto \prod_{i \in V} \phi^E(x_i) \prod_{(i,j) \in E} \phi(x_i, x_j)$$

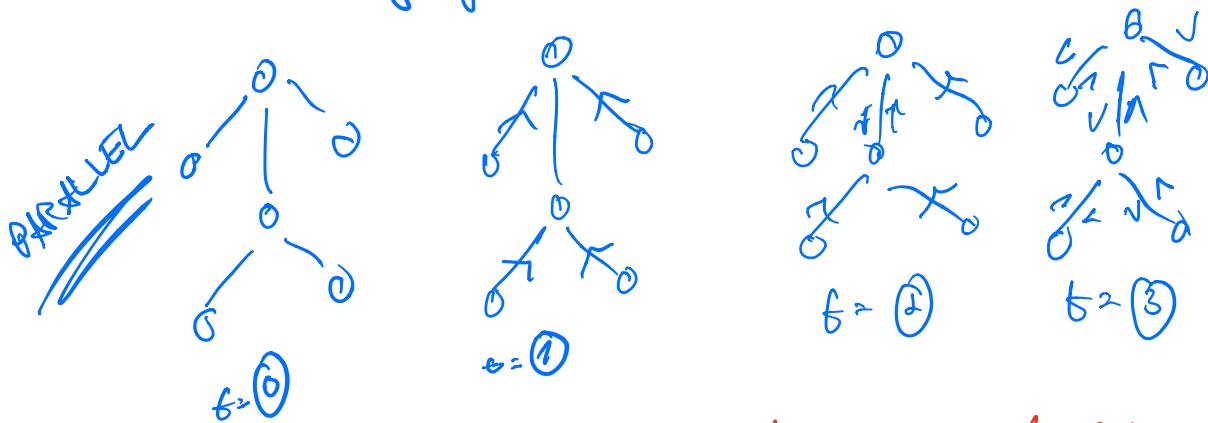
product of local factors

Message passing algo:

$$m_{j \rightarrow i}(x_j) = \sum_{x_j} (\phi^E(x_j) \phi(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}(x_j))$$

messages of other digraphs

Node can send message to neighbour after it has received messages from all of its remaining neighbors.



↳ Intract for non-tree structured graphs \Rightarrow Loopy
belief propagation \rightarrow variants: tree-re-weighted
 belief prop. / convex belief prop.

- MIP \rightarrow similar to marginal computation via Elimination
 ↳ maximize instead of sum // integrating

$$\mu_{ji}^{\text{mark}}(x_i) = \max_{x_j} \left(\phi^E(x_j) p(x_i, x_j) \prod_{k \in N(i)} \mu_{kj}(x_k) \right)$$

$$\max_x p^E(x) = \max_{x_i} \left(\phi^E(x_i) \prod_{j \in N(i)} \mu_{ji}^{\text{mark}}(x_i) \right)$$

- Variational Inference \leftrightarrow mean field / structured MF
 - \rightarrow convex opt. that is intractable since there are exp. many constraints, \Rightarrow connection of local to global
 - \rightarrow connection to sum-product

|o Take home: Exact inference is intractable!

- Message passing algs are very scalable since we perform many locally simple operations!

Talk 8 - P. Richterik - Random Optimization 30/08

- Empirical Risk Minimization:

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \text{loss}(l_x(a_i), b_i) + g(x)$$

↳ Monte Carlo interpretation of risk:

$$\min_{x \in \mathbb{R}^d} \mathbb{E}_{(a_i, b_i) \sim D} [\text{loss}(l_x(a_i), b_i)]$$

- Linear System Problem Formulation:

$$\underline{\text{LOSS}} \quad \mathbb{1}_{\{b_i\}}(t) = \begin{cases} 0, & t = b_i \\ +\infty \text{ otherwise.} \end{cases}$$

$$\underline{\text{REG}} \quad \frac{1}{2} \|x - x^*\|_B^2$$

- Legendre-Fenchel Transform \Leftrightarrow Convex Conjugate \rightarrow Dual
 \rightarrow Map from fct. to another fct.

↳ **Continued notes in the slides**

Folie 10 - S. Russell - Open-Universe prob. models 31/08

- first-order logic \Rightarrow fails too unreliable // memory to learn! like using propositional logic - circuits / GPUs
- probabilistic folks \Rightarrow Bayesian nets (BNT)
- First order prob. logicians \Rightarrow combine Bayesian nets with logical deductive

World has kings in it all we don't know what they are \Rightarrow 1st order open universe prob. logicians

↳ closed-universe: unique names, domain closed \Rightarrow PROLOG
↳ open-universe: allow uncertainty over existence \Rightarrow BLOG

Bayesian Logic

Talk 11 - C. Szepesvári - Bandits I

81/08/17

① Stochastic Finite-armed bandits

- Hosteller, Bush (1953) - animal learning (Mouse in T-maze)
- Thompson (1933) - Medical Trial // A/B Testing
→ don't throw away information → sequential design
- Lai, Robbins (1985) - UCB \Rightarrow avg. regret optimality
 \hookrightarrow fu et al (2002) - finite iter. result

- Bandits \Leftrightarrow Simplest version of problem of *!
 \Rightarrow decision making under uncertainty $\rightarrow \mathbb{E}$ vs. E
- Policy: Mapping of all possible histories to future actions! $\boxed{\text{algo}}$
- Min regret: O has special property \Rightarrow scaling invariant!
 \hookrightarrow still: scaling invariant
- O -regret lower: $\frac{R_n}{n} \rightarrow 0 \quad (\Rightarrow R_n = O(n))$

$\boxed{\text{speed}}$ \rightarrow How slow does R_n grow?

- Time horizon is important: No factor? Explanation does not make sense!

$A_t \sim \pi(\cdot | A_1, X_1, \dots, A_{t-1}, X_{t-1}) \Rightarrow$ agent action follows policy given history!
 $\Rightarrow r$: environment $\Rightarrow (r, \pi)$ joint distr. of interest!

- regret decomposition:
 - $\rightarrow V = (p_1, \dots, p_K)$: benefit environment → non-random
 - $\rightarrow \Delta_i(r) = \mu^*(r) - \mu_i(r)$: sub-opt. gap environ i
 - $\hookrightarrow T_i(t) = \sum_{s=1}^t \mathbb{1}_{\{\sum A_t = i\}}$ → random due to reward

\Rightarrow fixed π, V , then regret R_n :

$$R_n = \sum_{i=1}^K \Delta_i \mathbb{E}[T_i(n)] \quad \begin{matrix} \nearrow \\ \text{don't observe} \\ \text{behavior to estimate!} \end{matrix}$$

- Concentration of measure

$$\rightarrow \text{Markov: } P(|X| > \varepsilon) \leq \frac{\mathbb{E}[|X|]}{\varepsilon}$$

$$\rightarrow \text{Chebyshev: } P(|X - \mathbb{E}[X]| > \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

\rightarrow Subgaussian: X is σ -subgaussian if $\forall \lambda \in \mathbb{R}$:

$$\mathbb{E}[\exp(\lambda X)] \leq \exp(\lambda^2 \sigma^2 / 2)$$

\hookrightarrow factors are decaying factor!

\hookrightarrow If X

$$P(X \geq \varepsilon) \leq \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right)$$

\hookrightarrow fact.: R.V.s in bandit problem are sub-gaussian!

from
chebyshev
inequality.

①.1 Explore - flu - Commit \Rightarrow follow the leader after some exploration of all forms
 \Rightarrow ETC

- ϵ -greedy: part of the journey!
- visit number of iter. we are exploring

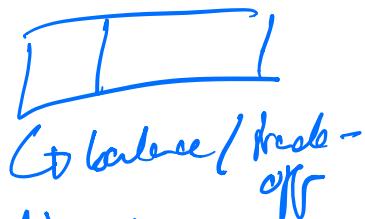
①.2 Justice - Adaptive Aggs

②. 1 UCB \Rightarrow Optimism in the face of uncertainty
 \rightarrow Choose actions as if environment is as nice as possible
 \rightarrow In beginning: each arm is explored once!

Talál 1d - C. Szepesvári - Backoff II Oktatás

- Recalling there: For each $L, \exists C > 0$ s.t.

$$R_n \leq C \sqrt{u}$$



- UCB \Rightarrow many different versions depending on how we accumulate data!

↳ simulation: balancing of upper bounds $\alpha \log(\text{Conf})$

$$UCB_i(t-1, \delta) = \hat{\mu}_i(t-1)$$

$$t \sqrt{2 \log(1/\delta)}$$

$$A_t = \arg\max_i \pi_{\text{SPLT}}^i$$

→ bad performance if environment is non-stationary
↳ bad can modify if we know the speed of change!

• Asympt. Conv.

• Finite time

• Worst Case



Tend to result in good algs!

→ Dangerous comment: Might be back to "only" minimize based on expected regret! \Rightarrow variance matters!

② Stochastic Bandits

→ Adversarial Environment:
 \forall point in $([0,1]^K)^n$: $\gamma = (\gamma_1, \dots, \gamma_n)$ where
 $x_t \in [0,1]^K$

↳ $t=0$: \forall chooses x_1, \dots, x_n

$t=1, 2, \dots$:

* Based on history, choose action $a_t \in [K]$

↳ send to \forall

* Env. sends reward $r_t = r_{x_t, a_t}$ to learner

\Rightarrow action chosen to minimize the worst-case!

$$R_n = \max_i \sum_{t=1}^n X_{t,i} - \mathbb{E} \left[\sum_{t=1}^n X_{t,i} \right]$$

choosing
 adversarial vs.
 choosing actions
 of
 non-random

rewards
 ...
 A

Important! \Rightarrow otherwise adversarial would simply
 exploit your deterministic behaviour! \Rightarrow actually
 have to randomize!

↳ reactive approach \rightarrow MDP formulation

$$R_n^*(\pi) = \sup_{\text{rewards}} \mathbb{E}_n(\pi, V) \rightarrow \text{worst-case regret}$$

\Rightarrow Existence: $R_n^*(\pi) = o(n) \rightarrow$ sublinear regret!

$$\mathbb{E} \left[\max_i \sum_{t=1}^n X_{t,i} - X_{t,\pi_t} \right] \geq \max_i \mathbb{E} \left[\sum_{t=1}^n X_{t,i} - X_{t,\pi_t} \right]$$

Exp. Adversarial regret
 \hookrightarrow It does not really matter if world is truly stochastic!

- Exp 3: Exponential-weight algo for Exploration & Exploitation
 \hookrightarrow Estimation of rewards \Rightarrow Only $X_t = X_{t,\pi_t}$ is seen
- Minimax regret: $R_n^*(\pi) = \inf_{\text{estimator}} \sup_{\pi \in \mathcal{E}} R_n(\pi, V)$

Talle 1B - S. Nowozin - Probabilistic DZ 01/08/17

- Evolution was not able / did not get enough time
 \Rightarrow It is all about representations and adaptive learning
- Lower dimensional subspaces capture abstract structures!
 \rightarrow accessible and controllable representations
 ↳ meaning: word2vec e.g.
 \rightarrow weak / implicit supervision
 \rightarrow robust learning
- tolerances in Density Est.

① Integral Prob. Metrics

$$Y_F(P, Q) = \sup_{f \in \mathcal{F}} |\int f dP - \int f dQ| \quad \text{→ Wasserstein GANs}$$

\rightarrow depending on class of function you get different metric

② proper Scoring Rules

$$SC(P, Q) = \int S(P, X) dQ(X)$$

\rightarrow most used: log likelihood!

→ ver. interval scales

③ f-divergences

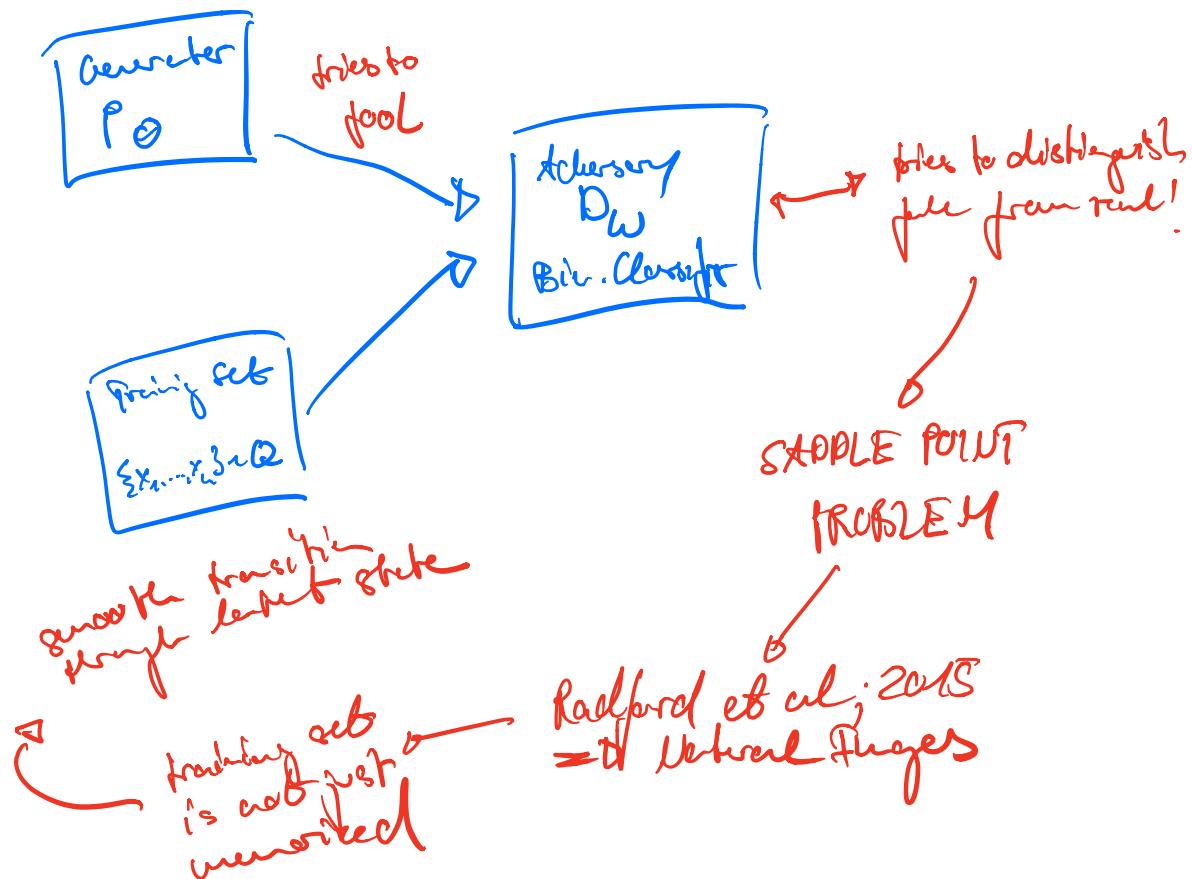
$$D_f(P||Q) = \int q(x) \underbrace{f\left(\frac{p(x)}{q(x)}\right)}_{\geq 0} dx \quad \text{→ OTUs}$$

\rightarrow KL divergence belongs to this class

f: generator
f: function
 $f(1) = 0$

* convex

- Implicit Model \Rightarrow fixed prob. dist. \rightarrow apply fct. to obtain a sample \rightarrow density fct.-weight not be defined explicitly
 - \hookrightarrow Kernel sampler / L1-free models $\xrightarrow{f(\cdot)}$
 - \rightarrow building block \Rightarrow very flexible
 - \rightarrow Problem: $p(x)$ might not be defined
 \hookrightarrow can't reason about all events \Rightarrow Bayes
not all dimensions can be recovered
 - \hookrightarrow possible solution: told value \rightarrow fct: how to choose?
- GAN = Implicit model \oplus Estimation procedure



- Relation to f-divergences
 \rightarrow Every convex fct. has a Fenchel conjugate!

$$P_f(P||Q) = \int_X q(x) \sup_{t \in \text{dom } f} \left\{ t \frac{p(x)}{q(x)} - f^*(t) \right\} dx$$

$$\geq \sup_{T \in \mathcal{T}} \left(\int_X p(x) T(x) dx - \int_X q(x) f^*(T(x)) dx \right)$$

$$= \sup_{T \in \mathcal{T}} \left(\mathbb{E}_{x \sim p}[T(x)] - \mathbb{E}_{x \sim Q}[f^*(T(x))] \right)$$

to approx. by MC integration / samples

$\Rightarrow f\text{-GAN} \rightarrow$ Jensen-Shannon Divergence:

$$\min_w \max_{\theta} \mathbb{E}_{x \sim p_\theta} [\log D_w(x)] + \mathbb{E}_{x \sim Q} [\log(1 - D_w(x))]$$

$\Rightarrow f\text{-GAN}$

$$\min_w \max_{\theta} \left(\mathbb{E}_{x \sim p_\theta}[T_w(x)] - \mathbb{E}_{x \sim Q}[f^*(T_w(x))] \right)$$

- Why does it work if $p(x)$ is not defined?
→ generalize f-clust. by introducing geometry of space
to Gaussian kernel \Rightarrow extended \rightarrow add noise
- GAN training is not always stable
→ numerical issues for ∇ player non-smooth
grads \Rightarrow solution: regularizers!
- Variational Autoencoders \rightarrow more stable \Rightarrow deleted to
derivative AAEs

$$\overline{p(x|\theta)} = \int p(x|z, \theta) \overline{p(z)/dz} \rightarrow \text{intractable integral}$$

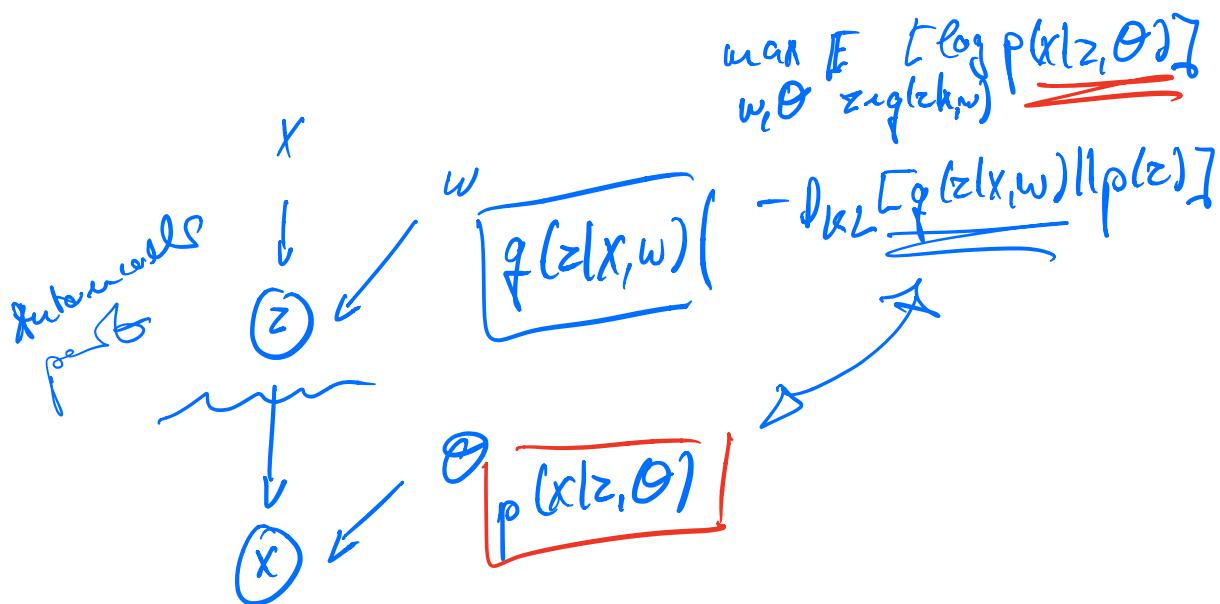
\hookrightarrow Mixture of Normals \hookrightarrow Multivariate Normal
 \Rightarrow fixed prior

\rightarrow MLE training \Rightarrow VI formulation

$$\begin{aligned}
 & \cancel{\log \int p(x|z, \theta) \frac{p(z)}{q(z)} q(z) dz} \\
 &= \log \mathbb{E}_{z \sim q(z)} \left[p(x|z, \theta) \frac{p(z)}{q(z)} \right] \\
 &\geq \mathbb{E}_{z \sim q(z)} \left[\log p(x|z, \theta) \frac{p(z)}{q(z)} \right] \quad \text{by Jensen's Ineq.} \\
 &= \mathbb{E}_{z \sim q(z)} \left[\log p(x|z, \theta) \right] - D_{KL}(q(z) || p(z))
 \end{aligned}$$

$\underbrace{\mathbb{E}_{z \sim q(z)} [\log p(x|z, \theta)]}_{\text{Exp. Var. free Energy}} \rightarrow$ Use surplus in practice!

→ Inference Network \Rightarrow new network ties together all conditional parameters \Rightarrow no direct optimization \rightarrow network does the optimization
 ↳ "parametrization"



→ Variable redefinition \Rightarrow reparameterization trick
 ↳ Doersch: "Tutorial on VAE"
 ↳ makes the gradients flow more smoothly!

- Problems:
 - * ELBO can be very loose!
 - * Image generation of VAEs \Rightarrow can be very blurry