

Francis Bach (INRIA) : Optimization

LARGE-SCALE
SUPERVISED ML

✓ ↗
LARGE n LARGE d

⇒ IDEAL: RUNNING -
TIME COMPLEX. $O(dn)$

$$\text{OPTIM.: } \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

↳ TRAINING ERROR

$$\text{STATISTICS: } E_{p(x,y)} l(y, h(x, \theta))$$

↳ TESTING COST

⇒ NEED TO RESPECT BOTH!

I. INTRO: SUPERVISED ML

□ DATA: $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}, i=1, \dots, n$
PRED. FCT.: $h(x, \theta) \in \mathbb{R}$ w. $\theta \in \mathbb{R}^d$

□ (REG.) EMPIRICAL RISK MINIMIZATION

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n l(y_i, h(x_i, \theta)) + \lambda \| \theta \|_2^2$$

$\rightarrow \theta^T \Phi(x)$

■ HINGE: $l(\cdot, \cdot) = \max\{1 - y h(x, \theta), 0\}$

■ LOGISTIC: $l(\cdot, \cdot) = \log(1 + \exp(-y h(x, \theta)))$

■ LS REG: $l(\cdot, \cdot) = \frac{1}{2} (y - h(x, \theta))^2$

↳ BIG ASS.: LINEARITY $h(\cdot)$ NOT CONVEX. LOSS

SMOOTHNESS + CONVEXITY + GRADIENT DESCENT

□ $g: \mathbb{R}^d \rightarrow \mathbb{R}$ IS L-SMOOTH IFF
ITS TWICE DIFF. + HESSIAN
IS BOUNDED:

$$\forall \theta \in \mathbb{R}^d, \text{ eigenv. } [g''(\theta)] \leq L$$

□ TWICE DIFF. $g: \mathbb{R}^d \rightarrow \mathbb{R}$ μ -STRONGLY
CONVEX IFF

$$\forall \theta \in \mathbb{R}^d, \text{ eigenv. } [g''(\theta)] \geq \mu$$

↳ CONDITION NUMBER DETERMINES
DIFFICULTY OF PROBLEM:

$$\kappa = L/\mu \geq 1$$

↳ CONVEXITY MAINLY ANALYSIS
TOOL → ALGOS SHOULD ALSO
RUN IN NON-CONVEX PROBLEM

↳ LINEAR CASE $[\Phi(x) \Phi(x)^T]^{-1} \Rightarrow n \gg d$

□ $g(\theta)$: OFTEN WEAKLY CONVEX

↳ STRONG-CONVEXITY RESULTS FROM
ADDING REGULARIZER!

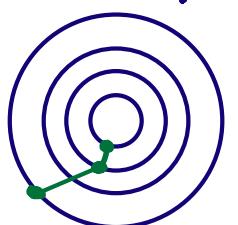
□ GRADIENT DESCENT:

↳ g CONVEX + L-SMOOTH ON \mathbb{R}^d

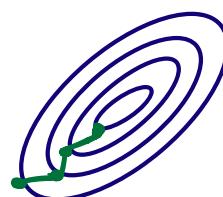
$$\theta_t = \theta_{t-1} - \gamma_t g'(\theta_{t-1})$$

LINE SEARCH $\rightarrow \frac{1}{n} \sum_{i=1}^n f_i'(\theta_{t-1})$

↳ FULLY ADAPTIVE!



SMALL κ



LARGE κ

⇒ LINEAR CONVERGENCE $\mathcal{O}(1/t)$

⇒ $\mathcal{O}(e^{-t/\kappa})$ IF STRONGLY CONVEX!

□ CHOOSE κ (#ITERS) PROP. TO CONDITION NUMBER κ !

□ ANALYZE $\|\theta_t - \theta_*\|^2$ VS. $g(\theta_t) - g(\theta_*)$

□ NEWTON 2ND ORDER METHOD: $\theta_t = \theta_{t-1} - g''(\theta_{t-1})^{-1} g'(\theta_{t-1})$

↳ "RECONDITION PROBLEM TO LOOK MORE LIKE BALL" \rightarrow NATURAL GRADIENTS

T-LGD	L-SMOOTH + μ -CONVEX	NON-STRONG CONV.
SGD	$d \times \frac{4}{\mu} \times \frac{1}{\epsilon}$	$d \times \frac{1}{\epsilon^2}$
BGD	$d \times n \frac{4}{\mu} \times \log \frac{1}{\epsilon}$	$d \times \frac{n}{\epsilon}$
ACC.GD	$d \times n \sqrt{\frac{4}{\mu}} \times \log \frac{1}{\epsilon}$	$d \times \frac{n}{\epsilon^2}$
SAG	$d \times (n + \frac{2}{\mu}) \times \log \frac{1}{\epsilon}$	$d \times \frac{\sqrt{n}}{\epsilon}$
ACC.	$d \times (n + \sqrt{n} \frac{4}{\mu}) \times \log \frac{1}{\epsilon}$	

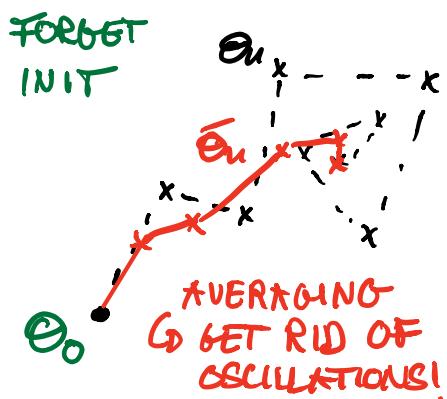
BOTTOU & BOUSQUET 08':
 * NO NEED TO OPTIM BELOW STAT. ERR
 * COST FUNCTIONS ARE AVERAGES
 * TESTING ERROR

II. FAST STOCHASTIC GRADIENT METHODS FOR CONVEX PROBLEMS

□ SGD \rightarrow POLYAK-RUPPERT AVG: $\bar{\theta}_t = \frac{1}{t+1} \sum_{u=0}^t \theta_u$

$$\hookrightarrow \mathbb{E} g(\bar{\theta}_t) - g(\theta_*) \leq \begin{cases} \mathcal{O}(1/t^2) & \text{if } \gamma_t = 1/(L\sqrt{t}) \\ \mathcal{O}(1/t) & \text{if } \gamma_t = 1/(\mu t) \end{cases}$$

MORE ITERS BUT SMALLER BATCH!



↳ RUNNING TIME: $\mathcal{O}(d \cdot \kappa / \epsilon)$

→ PROBLEM: NOT ADAPTIVE IN GENERAL!

□ GOAL: BEST OF STOCHASTIC + DETERM. WORLD

↳ LINEAR RATE w. $\mathcal{O}(d)$ ITER COST

⊕ SIMPLE CHOICE OF STEP CHOICE

□ ACCELERATION OF GRADIENT METHODS

\rightarrow NESTEROV 83', 04': $\theta_t = \gamma_{t-1} - \gamma_t g'(\gamma_{t-1})$ w. $\gamma_t = \theta_t + \delta_t^\top (\theta_t - \theta_{t-1})$

$$\begin{aligned} \hookrightarrow g(\theta_t) - g(\theta_*) &\leq \mathcal{O}(1/\epsilon^2) \\ g(\theta_t) - g(\theta_*) &\leq \mathcal{O}(e^{-\delta_t/\sqrt{K}}) \end{aligned}$$

STILL PROBLEM OF COMP. COMPLEXITY!

→ STOCHASTIC AVG. GRADIENT [LE ROUX ET AL., 12']

$$\Theta_t = \Theta_{t-1} - \frac{1}{n} \sum_{i=1}^n y_i^{(t)} w. y_i^{(t)} = \begin{cases} \hat{y}_i^{(t)}(\Theta_{t-1}) & \text{if } i = \text{idx} \\ y_i^{(t-1)} & \text{otherwise} \end{cases}$$

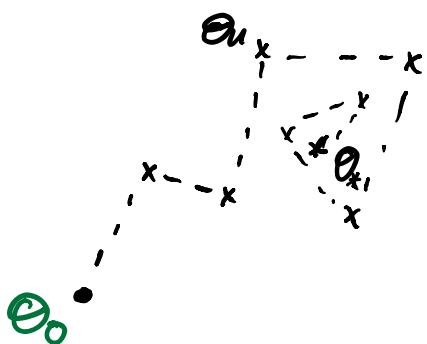
RANDOM SAMPLE
W. REPLACEMENT
ADD. MEMORY REQ.

↳ DIFFERENT INTERP.: LAZY GRAD EVAL. ⊕ VARIANCE REDUCTION

□ SGD MINIMIZES TEST COST! [INDEP. SAMPLES + SINGLE PASS \Rightarrow ROUND]

→ OPTIMAL CONV. RATE: $O(1/\sqrt{n})$ & $O(1/n\mu)$

→ CONSTANT STEP-SIZE SGD: LINEAR CONVERGENCE FOR STRONGLY CONV.



→ CONSTANT STEP-SIZE FOR LS:

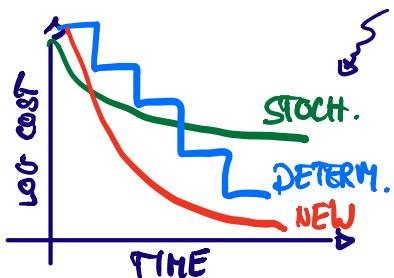
$(\Theta_n)_n \leftrightarrow$ HOMOGENOUS MARKOV CHAIN

↳ CONVERGENCE TO STATIONARY DISTR.

↳ ERGODIC THEOREM: AVG. ITERATES CONVERGE AT RATE $O(1/n)$

III. BEYOND CONVEX PROBLEMS

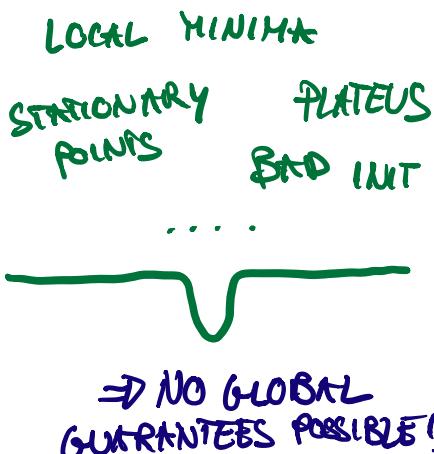
▷ SO FAR: CONVEX LOSS + LINEAR PREDICTOR $h(x, \Theta) = \Theta^T \Phi(x)$



↳ HOW TO TRANSFER SUCCESS TO DEEP LEARNING?

↳ NON-CONVEX OPTIMIZATION PROBLEM

↳ GENERALIZATION GUARANTEES IN OVERPARAMTERIZED REGIME?



□ LOCAL THEORETIC GUARANTEES

↳ REPLACE ITERATES BY FUNCTION VALUES!

↳ CHOROMSKA ET AL. '15: MOST LOCAL MINIMA ARE EQUIVALENT! \rightarrow MANY ASST.

↳ SOLTA-NOLKOTABLE ET AL. '18: NO SPURIOUS LOCAL MINIMA \rightarrow NOT USED ACTIVATIONS

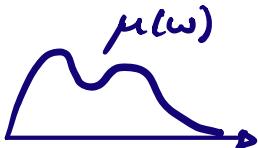
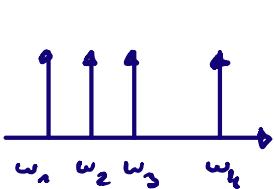
□ SINGLE HIDDEN LAYER FC NETWORK:

$$h(x) = \frac{1}{m} \Theta_2^T \sigma(\Theta_1^T x) = \frac{1}{m} \sum_{j=1}^m \Theta_2(j) \cdot \sigma[\Theta_1(:, j)^T x]$$

NO SHARING OF PARAMS

MEAN-FIELD LIMIT PERSPECTIVE!
[WEIGHTS = PARTICLES]

$$\sum_{j=1}^m \bar{\psi}_j(w_j) = \int_{\omega} \bar{\psi}(w) d\mu(w) \quad w \sim d\mu(w) = \frac{1}{m} \sum_{j=1}^m \delta_{w_j}$$



$$\min F(\mu) = R \left(\int_{\omega} \bar{\psi}(w) d\mu(w) \right)$$

↳ GRADIENT FLOW: $\dot{\mu} = -\eta \nabla F_m(w)$

↳ LIMIT $m \rightarrow \infty$: WASSERSTEIN GRADIENT FLOW

□ CHIZAT & BACH '18a: GRADIENT FLOW CONVERGES TO GLOBAL OPTIM.!

□ KEY ASSUMPTIONS: HOMOGENEITY [RELU/SIGMOID] + INIT MEASURE SPREAD
↳ NOTE: ONLY QUALITATIVE RESULT! → CAN WE GET MORE?

□ CHIZAT & BACH '18b: LAZY TRAINING [NOT TOO MUCH MOVEMENT F. INIT]

↳ SCALE FACTOR $\alpha > 0$: $G_\infty(\omega) = R(\alpha h(\omega)) / \alpha^2$

↳ INIT $w(0)$ s.t. $\alpha h(w(0))$ IS BOUNDED!

↳ AROUND $w(0)$, $G_\infty(\omega)$ has * ∇ "prop" to $\nabla R(\alpha h(w_0)) / \alpha$
* HESSIAN "prop" to $\nabla^2 R(\alpha h(w_0))$

⇒ MAKES MODEL LINEAR! → 1ST ORDER TAYLOR ARGUMENT

↳ NOT SO FAR FROM INIT!!!

□ JACOT ET AL. '18: NEURAL TANGENT KERNEL

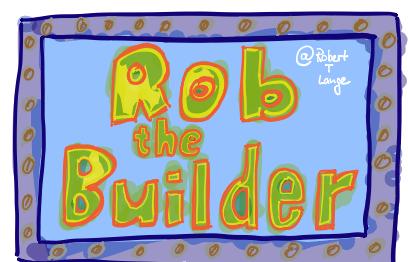
$$h(x, \omega) \approx h(x, w(0)) + [\omega - w(0)]^\top \nabla h(x, w(0))$$

$$\text{KERNEL: } k(x, x') = \nabla h(x, w(0))^\top \nabla h(x', w(0))$$

↳ EFFECT OF SCALE + APPLIES TO ALL DIFF. MODELS

□ PHYSICAL REVIEW: "NEW" / "FIRST" SHOULD NOT BE USED IN PAPERS

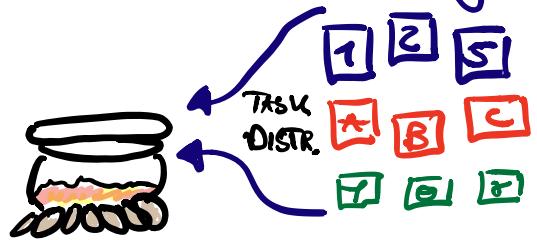
↳ DON'T OVERHYPE YOUR OWN WORK!



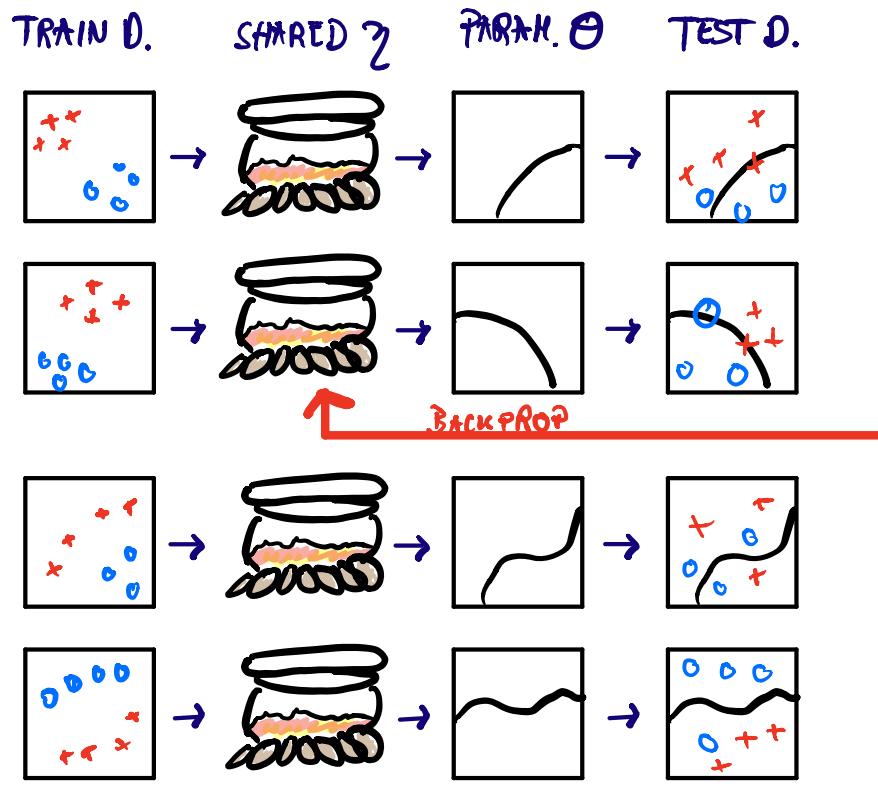
Yee Whye Teh (Ox+DeepMind) : Meta-learning

MOTIVATION: SMALL DATA PROBLEMS

- FEW-SHOT LEARNING
 - RECOMMENDER SYSTEMS
 - ROBOTICS
- } LEARN INDUCTIVE BIASES!



OPTIMISATION PERSPECTIVE ON META-LEARNING



INNER LEARN θ :

$$\begin{aligned}\theta_j &= \arg \min_{\theta_j} \sum_{i=1}^n \text{Loss}(f_{\gamma, \theta_j}(x_{ji}), y_{ji}) \\ &=: \text{Learner}(\gamma, \text{TrainData}_j)\end{aligned}$$

TEST PERFORMANCE:

$$\sum_{i \in \text{test}} \text{Loss}(f_{\gamma, \theta_j}(x_{ji}), y_{ji})$$

$$=: \alpha(\gamma, \theta_j, \text{TestData}_j)$$

OPTIM. SHARED PARAMS:

$$\begin{aligned}&\arg \min_{\gamma} \sum_{j=1}^T \alpha(\gamma, \theta_j, \text{TestData}_j) \\ &=: \arg \min_{\gamma} \sum_{j=1}^T \sum_{i \in \text{test}} \text{Loss}(f_{\gamma, \text{LEARNER}(\gamma, \text{Train})(x_{ji})}, y_{ji})\end{aligned}$$

OPTIMISATION-BASED META-LEARNING

□ BASE-LEARNER: θ_0

$$\theta_1 \leftarrow \theta_0 + \text{Update}_{\gamma} [\nabla_{\theta_0} L(.)]$$

$$\theta_2 \leftarrow \theta_1 + \text{Update}_{\gamma} [\nabla_{\theta_1} L(.)]$$

...

□ META-LEARN θ_0, γ

□ L+L GD-by-GD, HAML (Study choice 16', Fig 17')



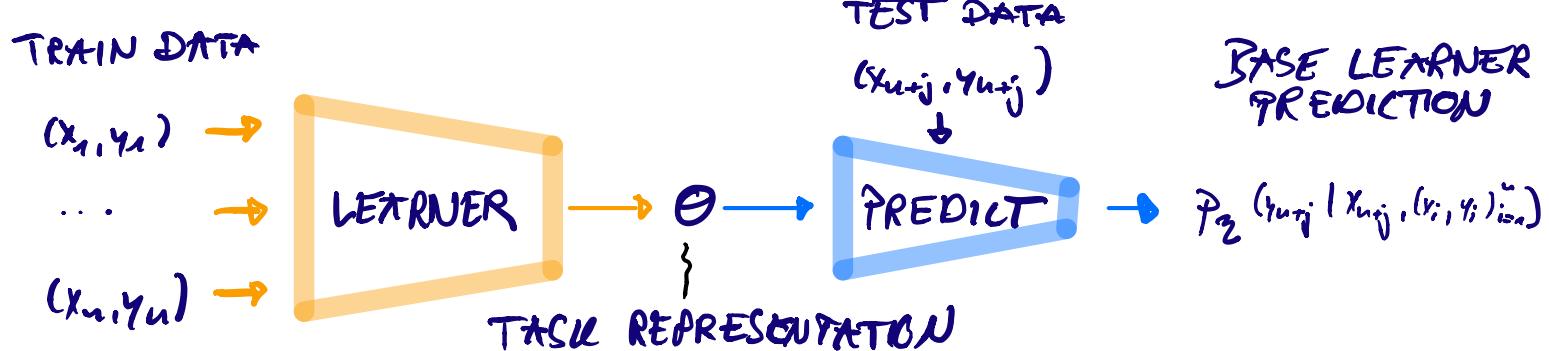
DIFFERENT INDUCTIVE BIASES

- PSYCHOLOGICAL THEORIES
- SYMMETRIES/INVARIANCES/EQUIVARIANCES
- COMP./LEARNTABILITY CONSTRAINTS

□ 2-LEVEL OPTIMISATION PROBLEM [INNER+OUTER]

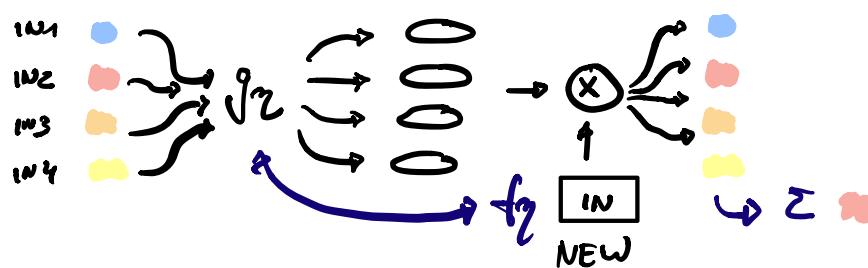
- ↳ FLEXIBLE + WIDELY APPLICABLE → CAN BE SENSITIVE TO ARCHITECTURE
- ↳ POSITIVE INDUCTIVE BIAS → CAN BE EXPENSIVE!

BLACK-BOX META-LEARNING → BASE-LEARNER = DIFFERENTIABLE PROGRAM



- ↳ REDUCES META-LEARNING TO SUPERVISED LEARNING ⇒ BROAD
- ↳ HARDER TO LEARN → NO INDUCTIVE BIAS ⇒ LESS GENERALIZATION

□ MATCHING NETWORKS [VINYALS ET AL. 16']



① EMBED TRAIN / TEST INPUT.

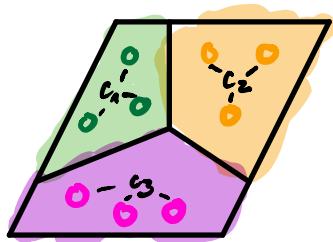
$$x_i \mapsto g_2(x_i)$$

$$x_{u+j} \mapsto f_2(x_{u+j})$$

② SOFT-1NN CLASSIFIER

$$\sum_{i=1}^n \frac{e^{g_2(x_i)^T f_2(x_{u+j})}}{\sum_{i=1}^n e^{g_2(x_i)^T f_2(x_{u+j})}} \tilde{y}_i$$

□ PROTOTYPICAL NETWORKS [SNELL ET AL. 17']



① EMBED $x_i \mapsto f_m(x_i)$

$$\text{② FORM PROTOTYPES } C_k = \sum_{i:y_i=k} f_m(x_i) / \sum_{i:y_i=k}$$

$$\text{③ PREDICT: } \frac{e^{-\|f_m(x_{u+j}) - c_k\|^2/\sigma^2}}{\sum_{k=1}^K e^{-\|f_m(x_{u+j}) - c_k\|^2/\sigma^2}}$$

□ MEMORY-AUGMENTED NNC [SANTORO ET AL. 16']

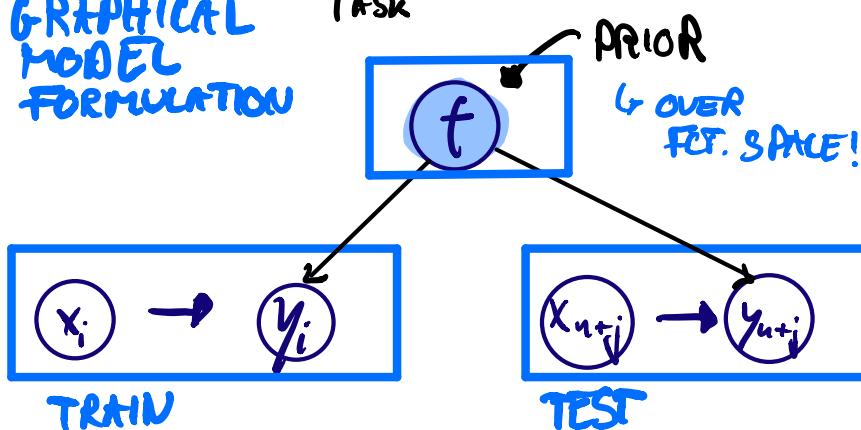
- TASK = SEQ. $(x_1, y_1), \dots, (x_T, y_T)$ ⇒ BASE: $p_{\pi_t}(y_t | (x_i, y_i)_{i=1}^{t-1}, x_t)$
- ↳ NEURAL TURING MACHINE WITH EXTERNAL MEMORY

□ CONDITIONAL NNs [GARREDO ET AL. 18'] → PERMUTATION INVARIANCE

- EMBED + AGGREGATE: $\Theta = 1/n \sum_{i=1}^n f_n(x_i, y_i) \leftrightarrow$ TASK REPR.

PROBABILISTIC PERSPECTIVE ON META-LEARNING

GRAPHICAL MODEL FORMULATION



\square STOCH. PROCESS = JOINT DISTR. OVER INF. COLLECTION OF RANDOM VARS. $[f(x), x \in \mathcal{X}]$

\square HOLMOGOROV EXTENSION TH.: CONSTRUCT STOCH. PROCESS BY SPECIFYING ITS FINITE DIM. MARGINAL DISTR.

\square FAMILY OF FINITE DIM. JOINT DISTR. $p_{x_{1:n}}$, ONE FOR EACH $n \in \mathbb{N}$ AND FINITE SEQ. $x_{1:n} \in \mathcal{X} \rightarrow$ WANT THESE TO FORM MARKONIALS:

$$p_{x_{1:n}}(y_{1:n}) = p(f(x_1) = y_1, \dots, f(x_n) = y_n)$$

\hookrightarrow EXCHANGEBILITY: FOR EACH n , $x_{1:n}$ AND PERMUTATION π OF $\{1, \dots, n\}$

$$p_{x_{1:n}}(y_{1:n}) = p_{x_{\pi(1)}, \dots, x_{\pi(n)}}(y_{\pi(1)}, \dots, y_{\pi(n)})$$

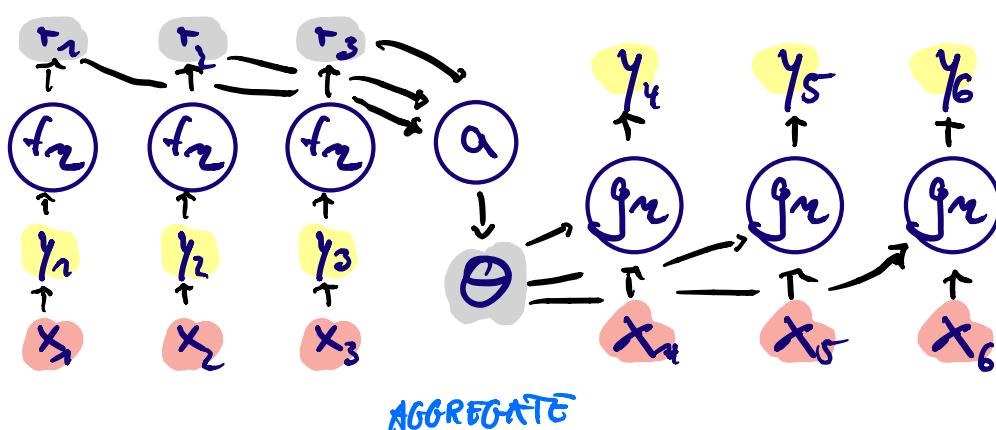
\hookrightarrow CONSISTENCY: FOR EACH n, m , $x_{1:n+m}$

$$p_{x_{1:n}}(y_{1:n}) = \int p_{x_{1:n+m}}(y_{1:n+m}) dy_{n+1:n+m}$$

\square GENERALLY: STOCH. PROCESS CAN BE DESCRIBED USING CONSISTENT FAMILY OF CONDITIONAL DISTRIBUTIONS:

$$p(f(x_{n+j}) = y_{n+j} | f(x_1) = y_1, \dots, f(x_n) = y_n)$$

\square CONDITIONAL NEURAL PROCESSES (II)



\rightarrow PREDICT

$$\begin{aligned} p_{\mathcal{Z}}(y_{n+j} | x_{n+j}, (x_i, y_i)_{i=1}^n) \\ = N(y_{n+j}; g_{\mathcal{Z}}(\Theta, x_{n+j})) \end{aligned}$$

\rightarrow PROBLEM: ARCH. CANNOT MODEL DEP. AMONG TEST OUTPUTS

$$\begin{aligned} p_{\mathcal{Z}}(z, y_{1:n+m} | x_{1:n+m}) &= p_{\mathcal{Z}}(z) \prod_{i=1}^{n+m} p_{\mathcal{Z}}(y_i | z, x_i) \quad \text{GEN. LATENT} \\ &\leftarrow \text{VAR. MODEL} \qquad \qquad \qquad \text{AUTOREGRESSIVE} \\ \hookrightarrow \text{ELBO ON } \log p(y_{n+1:n+m} | x_{1:n+m}, y_{1:n}) \end{aligned}$$

□ META-LEARNING AS LEARNING STOCH. PROCESS \rightarrow PRIOR OVER FCTS.

↳ THINK OF BASE LEARNER AS FORM OF AMORTIZED LEARNING

↳ IMPORTANCE OF UNCERTAINTY: ACTIVE LEARNING, BO, RL

□ EFFICIENT MODEL-BASED RL \Rightarrow META-LEARN $\hat{p}(s'|s, a), \hat{r}(s, a)$

↳ ADVERSARIAL TESTING OF RL AGENTS \rightarrow BAYES OPTIMIZATION PROBLEM \Rightarrow GALTSHOV ET AL. 19' \rightarrow MAZE CONSTRUCTION

PROBABILISTIC SYMMETRIES AND NEURAL ARCHITECTURES

□ CHARACTERISING PERMUTATION-INARIANT FCTS. $h: X^n \rightarrow Y$

\rightarrow IS PERM.-INV.: $h(\pi \cdot (x_1, \dots, x_n)) = h(x_1, \dots, x_n)$

\rightarrow IS PERM.-EQUIV.: $h(x_1, \dots, x_n) = (y_1, \dots, y_n) \rightarrow h(\pi \cdot (x_1, \dots, x_n)) = \pi \cdot h(x_1, \dots, x_n)$

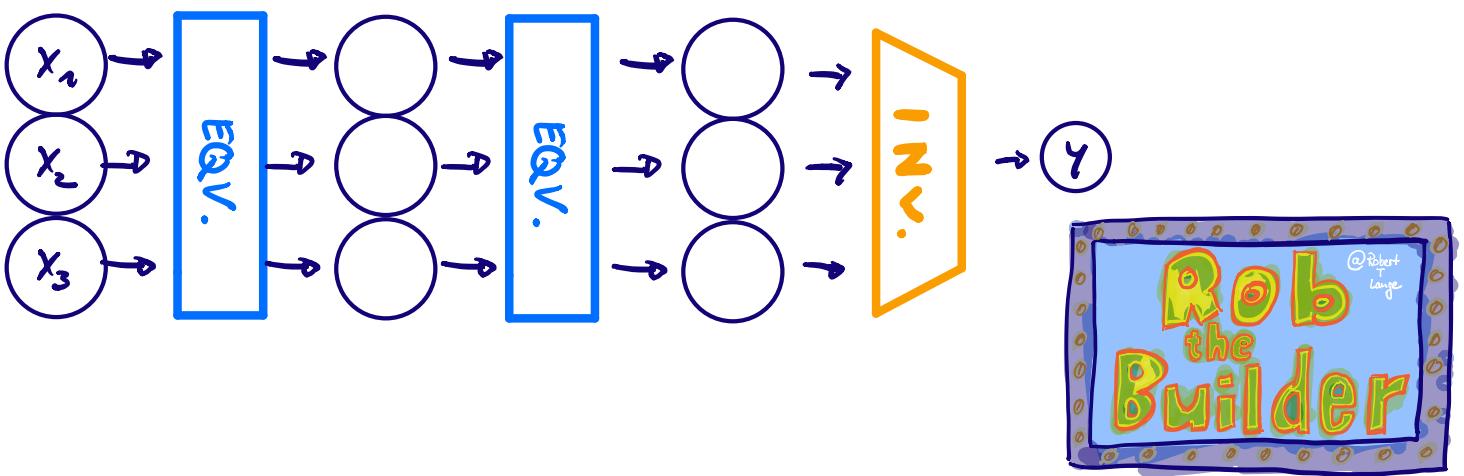
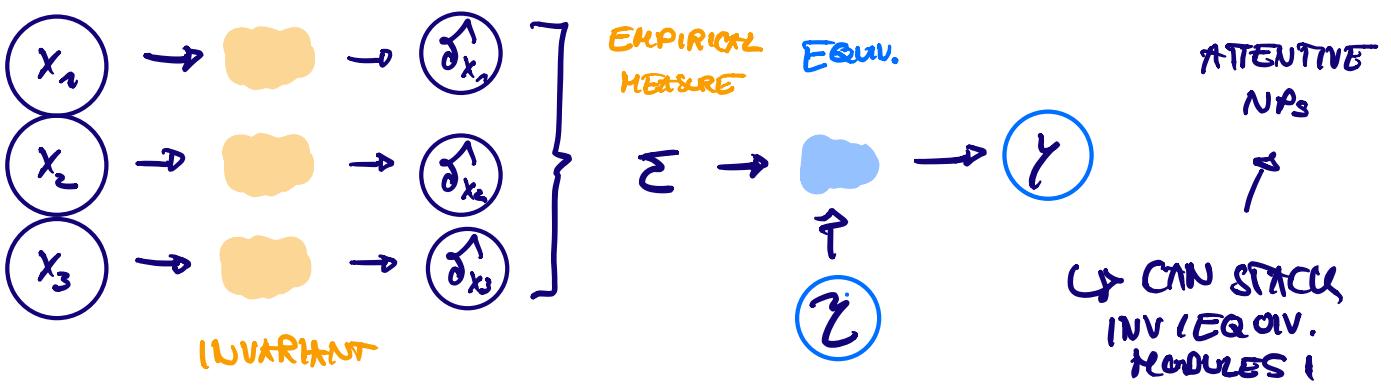
□ \neq COND. DISTR. $P(Y|X)$ IS STOCH. RELAX. FOR FCT. $Y = h(X)$

$\rightarrow P(Y|X)$ IS G-INVARIANT IF: $P(Y|X) = P(Y|g \cdot X)$

$\rightarrow P(Y|X)$ IS G-EQUIVARIANT IF: $P(Y|X) = P(g \cdot Y|g \cdot X)$

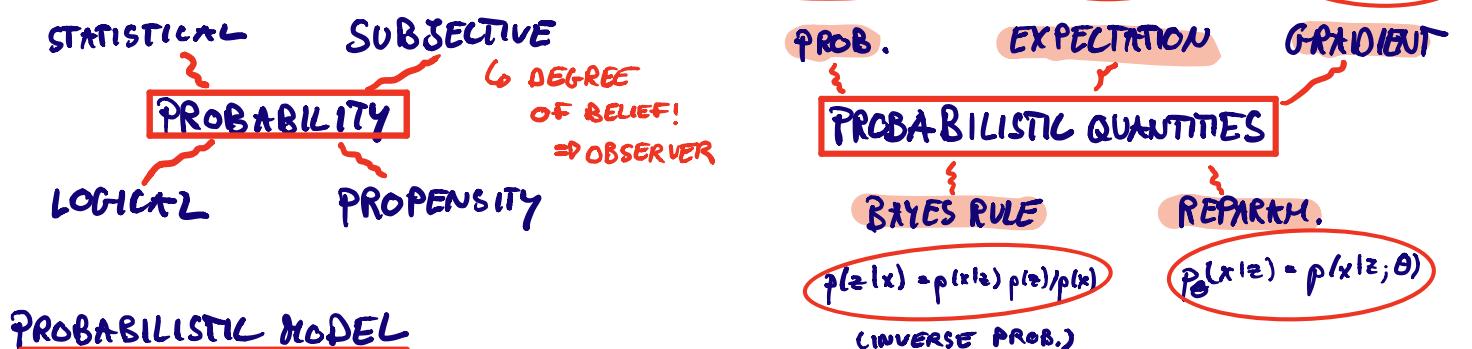
□ EMP. MEASURE OF X_n : $M_{X_n}(\cdot) = \sum_{i=1}^n \delta_{X_i}(\cdot)$

↳ SUFFICIENT STATS \rightarrow ADEQUATE STATS \rightarrow NOISE OUTSOURCING

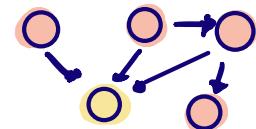


Shakir Mohamed (DeepMind): Bayesian Learning

I. BAYESIAN BASICS



→ MODEL = DESCRIPTION OF WORLD, DATA, POTENTIAL SCENARIOS, PROCESSES
 ↳ PROBABILISTIC: USES LANGUAGE OF PROB. TO WRITE OUT!
 ⇒ AIM = LEARN PROB. DISTR. OF DATA → REASONING



PROBABILITY OF A SEQUENCE



COIN Toss SEQUENCE

1, 0, 1, 1, ...

→ EXCHANGEABILITY: $p(x_1, \dots, x_n) = p(x_{\pi(1)}, \dots, x_{\pi(n)}) \Rightarrow$ INFINITE - JOIN! INV. TO PERMUTATION
 → DEFINITION: $p(x_1, \dots, x_n) = \int_{\Theta}^N p(x_n | \theta) P(d\theta)$ ⇒ PARAMETERS, PRIORS, AVERAGE

$$p(x_1, \dots, x_n) = \int_{\Theta}^N p(x_n | \theta) p(\theta) d\theta$$

ASSUMPTION OF DENSITY!

DATA x_1, \dots, x_n IS CONDITIONALLY INDEP. LIKELIHOOD $p(x | \theta)$ PRIOR $p(\theta)$ PARAMETER θ

MODEL-BASED!
JOINT PROB. VIA PARAM. LIKELIHOOD

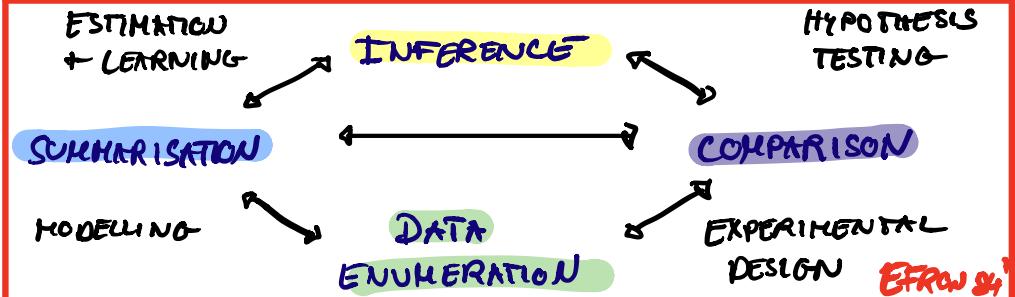
→ EVIDENCE: $p(y | x) = \int p(y | h(x); \theta) p(\theta) d\theta \rightarrow$ INTEGRATE OUT UNOBS.
 ↳ OFTEN INTRACTABLE

→ POSTERIOR: $p(\theta | y, x) \propto p(y | h(x); \theta) p(\theta)$

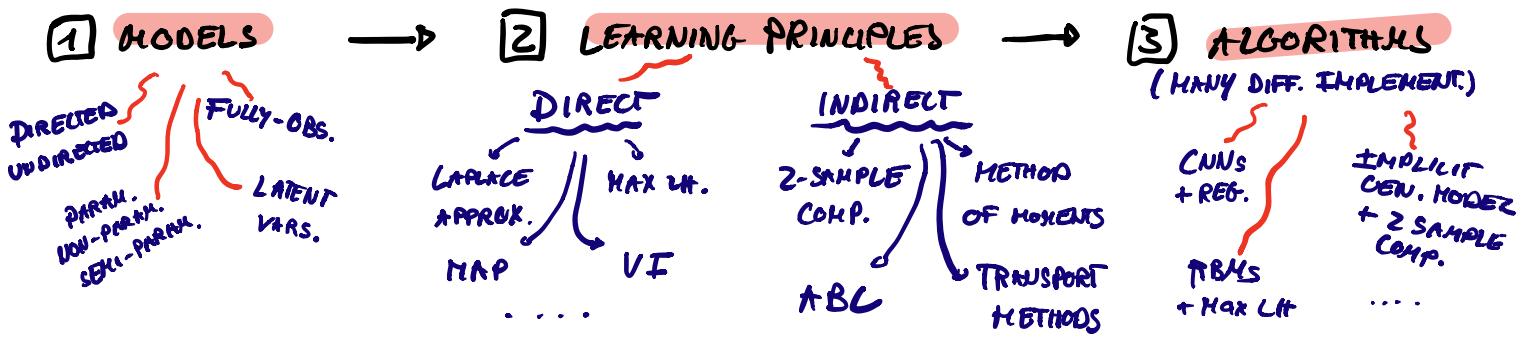


MARR'S LEVELS OF ANALYSIS:
COMPUTATIONAL ALGORITHMIC IMPLEMENTATION

SHAKIR'S 4 LEVELS:
APPLICATIONS
REASONING
INFORMATION PRINCIPLES

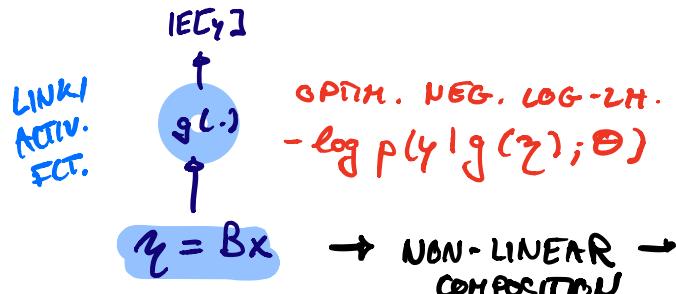


MODELS — LEARNING PRINCIPLES — ALGORITHMS

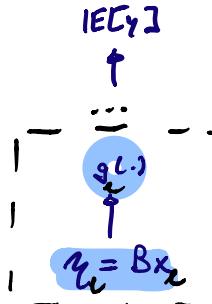


II. BAYESIAN COMPUTATION

LINEAR REGRESSION



DEEP LEARNING



HIERARCHICAL MODELS

GR DECOMP. OF PROB. DISTR.
INFO SEQ. OF CONDITIONAL DISTRIBUTIONS



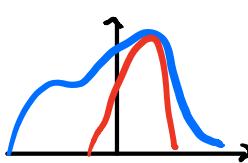
→ LIKELIHOOD FUNCTION: $\ell(\theta) = \sum_u \log p(y_u|x_u; \theta) \rightarrow$ LIKELIHOOD OF PARAMS.

- EFFICIENT ESTIMATORS
- TESTS W. GOOD POWER
- POOR INFORMATION
- WIDELY-APPPLICABLE
- MISSPECIFICATION POSSIBLE

↳ MAX LH: $\max_{\theta} \ell(\theta) \Rightarrow$ MAP: $\max_{\theta} \ell(\theta) + 1/2 R(\theta) \Rightarrow$ SHRINKAGE

→ INVARIANT MAP: $p(y|u(x); \theta) p(\theta) / I(\theta)^{1/2} \rightarrow$ REMOVE PARAM. SENSITIVITY

→ BAYESIAN INFERENCE = EVIDENCE $\Rightarrow p(x) = \int p(x|\theta) p(\theta) d\theta$ & POSTERIOR: $p(\theta|x)$



1 CONJUGACY = POSTERIOR IS IN SAME FORM OF OBS. DISTRIB.

↳ TYPICALLY: CONJUGATE PRIORS FALL INTO EXPONENTIAL DISTRIBS.

2 INTEGRAL APPROX. → E.G. 1ST ORDER TAYLOR \Rightarrow LAPLACE!

→ LEARNING AND INFERENCE IN MACHINE LEARNING

- INFERENCE = REASON ABOUT / COMPUTE UNKNOWN PROB. DISTR.
- LEARNING = FINDING GOOD POINT ESTIMATES OF MODEL QUALITIES

INFERENTIAL QRS!

EVIDENCE ESTIMATION

$$p(x) = \int p(x, z) dz$$

PREDICTION

$$p(x_{t+1}|x_{0:t})$$

MOMENT COMPUTATION

$$E[f(z)|x] = \int f(z) p(z|x) dz$$

PLANNING

$$IE_p[\int C(t_z) dt | x_{0:t}] \log p(x|t_z) - \log p(x|t_0)$$

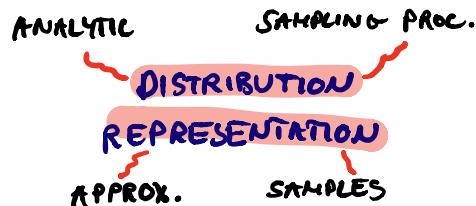
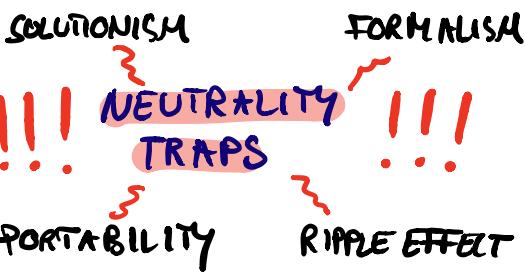
PARAMETER ESTIMATION

$$p(\theta|x_{0:N})$$

HYPOTHESIS T.

EXPERIMENTAL DESIGN

$$D[p(x_{t:T}|u) || p(x_{0:t})]$$



→ MONTE CARLO METHODS:

- ↳ INTEGRAL: $\hat{f}(\theta) = \int f(x) p(x|\theta) dx$
- ↳ SAMPLE ESTIM.: $\hat{f}(\theta) = \frac{1}{N} \sum_{i=1}^N f(\hat{x}^i), \hat{x}^i \sim p(x|\theta)$

DRAW FROM $p(x|\theta)$

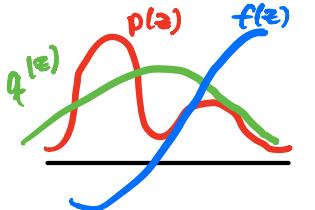
CONSISTENCY
ESTIMATE → TRUE

UNBIASEDNESS
TRUE CENTRING

LOW VARIANCE
ESTIM. ↔ RV

COMPUTATION
EFFICIENT

IMPORTANCE SAMPLING



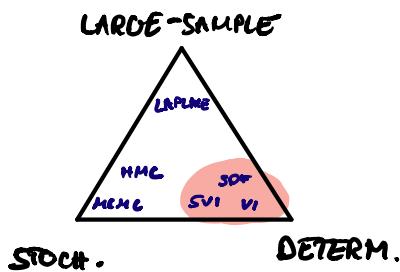
$$m = \mathbb{E}_{q(z)} \left[f(z) \frac{p(z)}{q(z)} \right]$$

↳ IDENTITY TRICK
⇒ EASY SAMPLING!

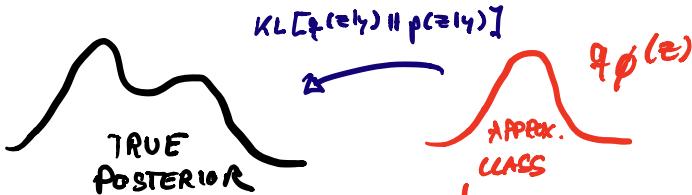
LIMITATIONS:
* COMP. INTENSIVE
* MIXING GUARANTEES

* SOLVE BAYES PROBLEM USING FREQU. METHODS // PAST DEP. UNKNOWN OR

↳ MANY DIFFERENT VERSIONS ⇒ QUASI-MC, BAYESIAN MC



VARIATIONAL METHODS



→ DETERM.
APPROX. PROCEDURE
↳ FOR PARAMS!

FROM IMPORTANCE SAMPL.
TO VARIATIONAL INF.

$$\log p(x) = \int p(x|z) p(z) / q(z) dz [IS]$$

$$\log p(x) \geq \int q(z) \log (p(x|z) p(z) / q(z)) dz [JENSEN'S INEQ.]$$

$$= \int q(z) \log p(x|z) - \int q(z) \log \frac{q(z)}{p(z)}$$

→ TURN INTEGRATION INTO OPTIM. PROBLEM

$f(x, q)$

⇒

$$\mathbb{E}_{q(z)} [\log p(x|z)] - KL[q(z) || p(z)]$$

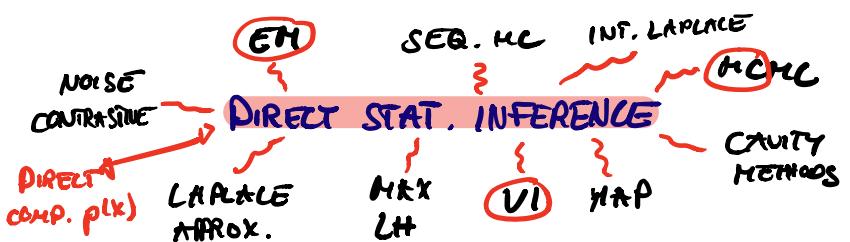
APPROX. POSTERIOR

RECONSTRUCTION

PENALTY → SMOOTH

[VARIATIONAL
LOWER BOUND]

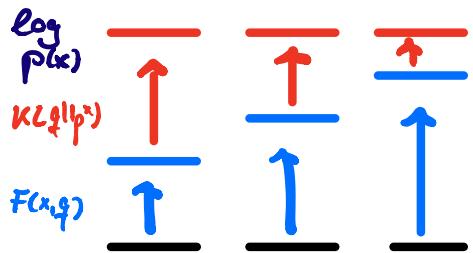
III. BAYESIAN APPROXIMATION



CLASSICAL INFERENCE APPROACH

$$p(x, z) \rightarrow \int (\dots) q_\phi(z|x) dx \quad \xrightarrow{\text{E-STEP}} \quad \nabla_\phi \quad \xrightarrow{\text{M-STEP}}$$

MODERN METHODS
↔ SVI

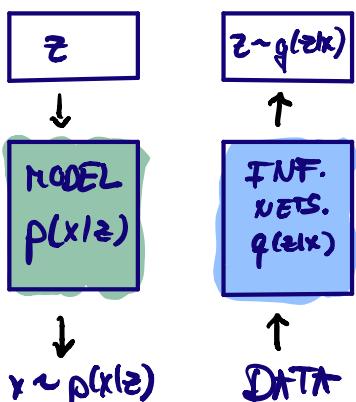


LOG-DERIVATIVE TRICK: $\nabla_\phi \log q_\phi(z) = \frac{\nabla_\phi q_\phi(z)}{q_\phi(z)}$

) REWRITE SO THAT
1. SAMPLE + 2. EVALUATE

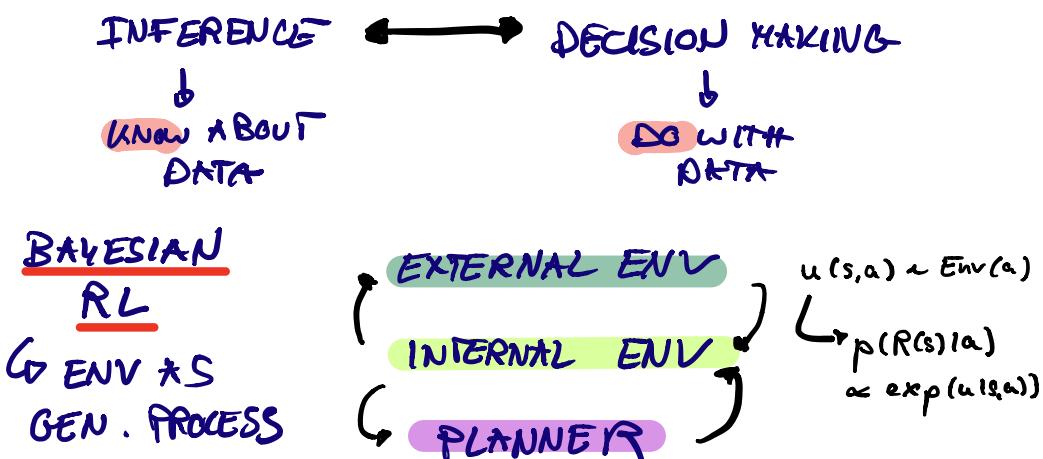
SCORE FCT. GRADIENT: $\nabla_\phi I\mathbb{E}_{q_\phi(z)} [f_\phi(z)] = I\mathbb{E}_{q_\phi(z)} [f(z) \nabla_\phi \log q_\phi(z)]$

VAEs



\Rightarrow SPEC. COMB. OF VI
IN LATENT VAR. MODEL
USING INF. NETS

IV. BAYESIAN FUTURES

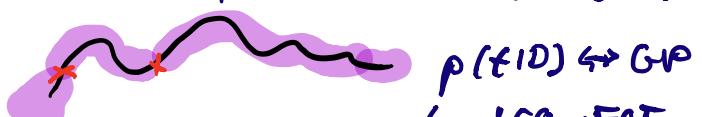


$$\mathcal{F}(\theta) = I\mathbb{E}_{\pi(a|s)} [R(s,a)] - KL[\pi_\theta(a|s) || p(a)]$$

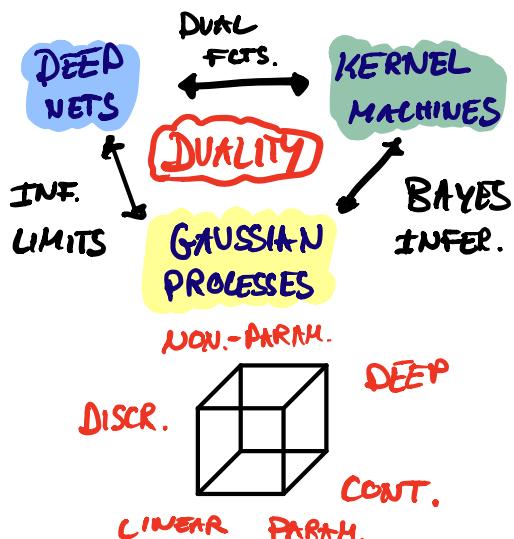
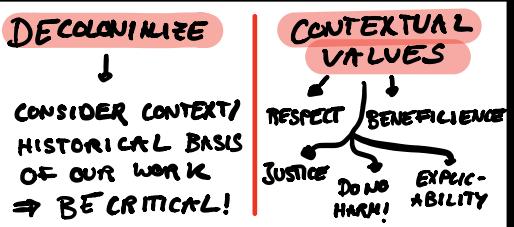
$$\nabla_\theta \mathcal{F}(\theta) = I\mathbb{E}_{\pi(a|s)} [(R(s,a) - \bar{c}) \nabla_\theta \log \pi_\theta(a|s)] + \nabla_\theta \pi_\theta(a|s)$$

\Rightarrow HIERARCHICAL PLANNING VIA ACTION PLANNING

BAYESIAN OPTIMIZATION \Rightarrow GLOBAL OPTIM. OF BLACK-BOX $f(x)$



$p(f|D) \leftrightarrow \text{GP}$
 \hookrightarrow ACQ. FCT.



PROB. NUMERICS

\Rightarrow THINK OF NUMERICAL METHODS AS BAYESIAN LEARNING ALGORITHMS

\hookrightarrow ODES / PDES / LINE SEARCH

