

Implicit Quantile Networks (IQN)

- Prev work (QR-DQN) learnt a distributⁿ over returns but still used mean returns for policy.
- IQN learns a risk sensitive policy.

Risk-Sensitive RL - I

- Instead of $\max E[\text{Returns}]$, $\max E[\text{Utility}]$.
- E.g.s of Utility: Exp/Log(Returns), Measure that includes variance etc.
- In general, utility f^n satisfies 4 axioms
 - ① Monotonicity: lower cost \Rightarrow lower risk
 - ② Sub-additivity: Diversification \Rightarrow lower risk
 - ...
- Can be concave(risk averse)/convex(seeking) or both.

Risk-Sensitive RL - II

- One more axiom of utility f^n is of independence:
if $X \succ Y$ (preferred over) Z , then $pX + (1-p)Z \succ pY + (1-p)Z$
(Any mixture of $X \& Z$ preferred over same mixture over $Y \& Z$).
We can think of utility as mixture of distributions, which we want to learn, to learn a risk-sensitive policy.

(As before, we want to use quantiles to learn distributⁿ)

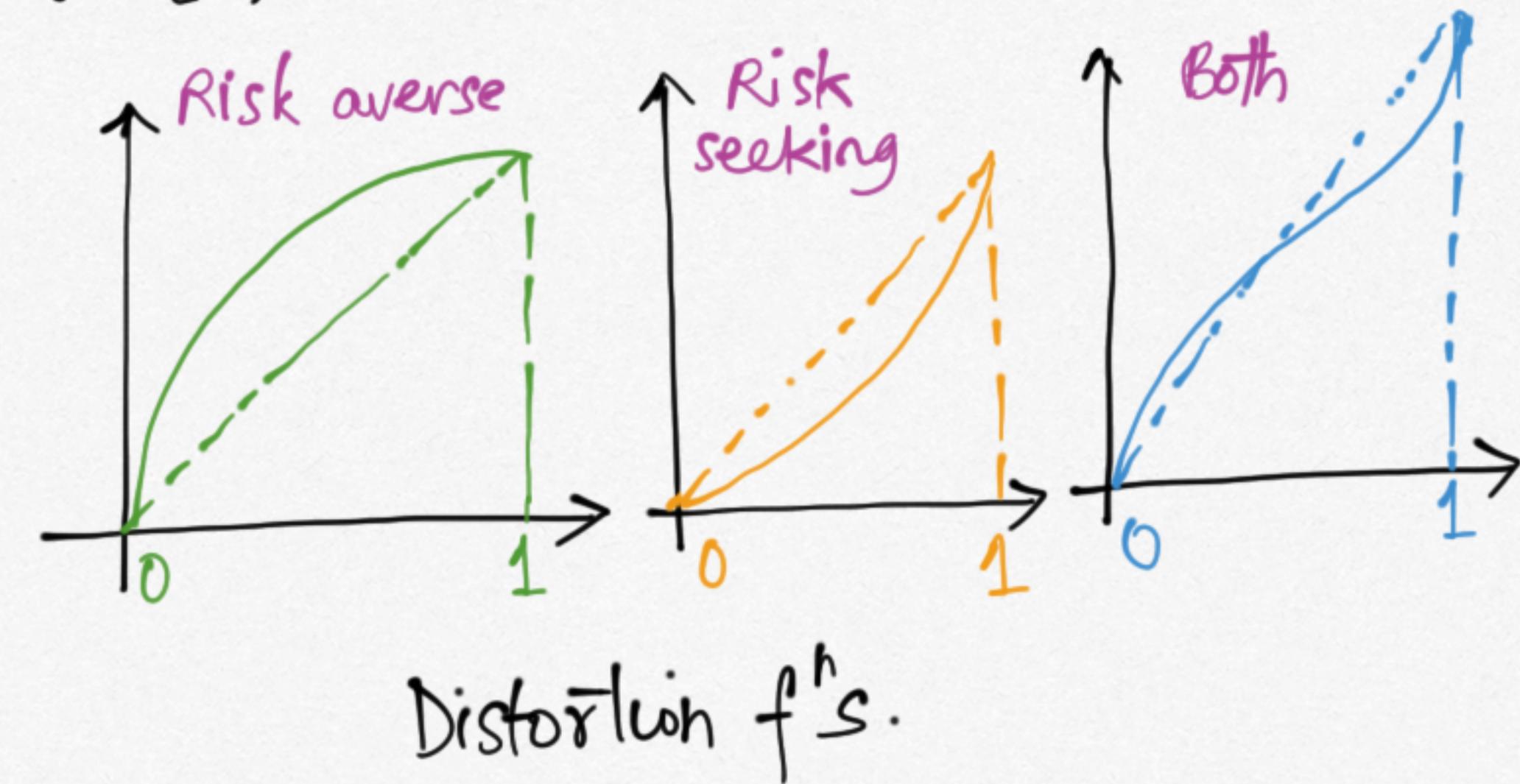
Learning mixture of distributions (MoD)

- MoD can be represented in terms of cdf.
 - cdf of MoD is convex combination of its individual distributions.
- $$F(x) = \sum_i p_i F_i(x)$$
- MoD is also represented as convex comb. of quantiles (in insurance/reinsurance...) where weighting is some risk fⁿ. (e.g. VaR)

Distribution Risk Measures (DRM)

- Var, Conditional Var are DRMs.
- In general, DRM is a continuous non-decreasing fⁿ.

$$g: [0,1] \rightarrow [0,1] \text{ s.t } g(0)=0, g(1)=1.$$



Dual theory of Choice under Risk

- An alternate objective for risk sensitive policy:
Instead of maximizing expected utility,
maximize expectation of distortion risk measure
 \Rightarrow i.e. Distorted Expectation.

Learning a Distorted Expectation

Say $Z(s, a)$ is our usual returns distribution.
for some quantile τ , $E[Z_\tau(s, a)]$ is expected return
 $\tau \in U(0, 1)$ for that quantile.

We can reweight this expectation w/ β (DRM),
 $Q_\beta(s, a) := E[z_{\beta(\tau)}(s, a)]$

[Can think of this as Inverse Transform Theorem:
To generate samples from some distribution, sample
from $Unif(0, 1)$, & transform them through inverse
CDF (quantile) of the prob. distribution.]

$Q_\beta(s, a)$ can be used to follow policy (action-selection)
 $\pi_\beta(s) = \operatorname{argmax}_{a \in A} Q_\beta(s, a)$

IQN Loss function

- Similar to QR-DQN we want to minimize loss b/w pairs of quantiles. But we don't have a fixed set of quantiles as in QR-DQN
- Instead, parametric $f^n \phi(t)$ used to represent quantiles (embedding layer in NN).
- So we can generate samples, N for τ , N' for τ' & calc. loss.
$$L = \frac{1}{N'} \sum_{i=1}^N \sum_{j=1}^{N'} p_{\tau_i}^k (\delta^{\tau_i, \tau'_j})$$

 $\delta^{\tau, \tau'} - TD$ error

- Similarly to select actions, use K samples of $\tilde{\tau}$ & take argmax of mean $Q_{\beta(\tilde{\tau})}(s, a)$
- How to generate these samples?
Sample from $U(0, 1)$ & reweight using the parametric f .