# Bandit Structured Prediction for Learning from Partial Feedback in Statistical Machine Translation

**Changzhi Sun**

- Expected Loss Minimization under Full Information

$$\mathbb{E}_{p(x,y)p_w(y'|x)}\left[\Delta_y(y')\right] = \sum_{x,y} p(x,y) \sum_{y' \in \mathcal{Y}(x)} \Delta_y(y')p_w(y'|x)$$

- minimum Bayes risk principle

$$\hat{y}_w(x) = \arg\min_{y \in \mathcal{Y}(x)} \sum_{y' \in \mathcal{Y}(x)} \Delta_y(y')p_w(y'|x).$$

- if $\Delta_y(y') = \mathbf{1}[y \neq y']$, MAP

$$\hat{y}_w(x) = \arg\max_{y \in \mathcal{Y}(x)} p_w(y|x)$$
$$= \arg\max_{y \in \mathcal{Y}(x)} w^\top \phi(x,y).$$

- $p(x, y)$ is approximated by
$$\tilde{p}(x, y) = \frac{1}{T} \sum_{t=0}^{T} \mathbf{1}[x = x_t]\mathbf{1}[y = y_t]$$
- training data $\{(x_t, y_t)\}_{t=0}^{T}$

$$\mathbb{E}_{\tilde{p}(x,y)p_w(y'|x)}\left[\Delta_y(y')\right] = \frac{1}{T} \sum_{t=0}^{T} \sum_{y' \in \mathcal{Y}(x_t)} \Delta_{y_t}(y')p_w(y'|x_t).$$

$$\nabla \mathbb{E}_{\tilde{p}(x,y)p_w(y'|x)}\left[\Delta_y(y')\right]$$
$$= \mathbb{E}_{\tilde{p}(x,y)}\left[\mathbb{E}_{p_w(y'|x)}[\Delta_y(y')\phi(x,y')] - \mathbb{E}_{p_w(y'|x)}[\Delta_y(y')]\,\mathbb{E}_{p_w(y'|x)}[\phi(x,y')]\right]$$
$$= \mathbb{E}_{\tilde{p}(x,y)p_w(y'|x)}\left[\Delta_y(y')(\phi(x,y') - \mathbb{E}_{p_w(y'|x)}[\phi(x,y')])\right].$$

- Bandit Structured Prediction

---

**Algorithm 1** Bandit Structured Prediction

1: Input: sequence of learning rates $\gamma_t$
2: Initialize $w_0$
3: **for** $t = 0, \ldots, T$ **do**
4:     Observe $x_t$
5:     Calculate $\mathbb{E}_{p_{w_t}(y'|x_t)}[\phi(x_t, y')]$
6:     Sample $\tilde{y}_t \sim p_{w_t}(y'|x_t)$
7:     Obtain feedback $\Delta(\tilde{y}_t)$
8:     Update $w_{t+1} = w_t - \gamma_t \, \Delta(\tilde{y}_t)(\phi(x_t, \tilde{y}_t) - \mathbb{E}_{p_{w_t}(y'|x_t)}[\phi(x_t, y')])$

---

$$J(w) = \mathbb{E}_{p(x)p_w(y'|x)}\left[\Delta(y')\right]$$
$$= \sum_x p(x) \sum_{y' \in \mathcal{Y}(x)} \Delta(y')p_w(y'|x).$$

- Structured Dueling Bandits

---

**Algorithm** Structured Dueling Bandits

---

1: Input: $\gamma, \delta, w_0$
2: **for** $t = 0, \ldots, T$ **do**
3:     Observe $x_t$
4:     Sample unit vector $u_t$ uniformly
5:     Set $w'_t = w_t + \delta u_t$
6:     Compare $\Delta(\hat{y}_{w_t}(x_t))$ to $\Delta(\hat{y}_{w'_t}(x_t))$
7:     **if** $w'_t$ wins **then**
8:         $w_{t+1} = w_t + \gamma u_t$
9:     **else**
10:         $w_{t+1} = w_t$

---

# Discussion

- first proposed **Bandit Structured Prediction**
- **single** point feedback