

Goal

We want to make **efficient** online learning algorithms that **adapt** automatically to the complexity of the environment.

- Worst-case rates in adversarial environments (safe and robust)
- Fast rates in favorable stochastic environments (practice)

Key Observation

Modern adaptive algorithms bound **regret** in terms of **variance**



Fast Rates

Friendliness of stochastic environments commonly quantified by condition relating **variance** to **regret**

Learning Model: Online Convex Optimization

In round $t = 1, 2, \dots$

- Learner chooses $w_t \in \mathcal{U} \subseteq \mathbb{R}^d$
- Environment selects convex loss function $\ell_t : \mathcal{U} \rightarrow \mathbb{R}$
- Learner incurs loss $\ell_t(w_t)$ and observes gradient $\nabla \ell_t(w_t)$

Goal: small **regret** R_T^u (or upper bound \tilde{R}_T^u) w.r.t. **every point** u

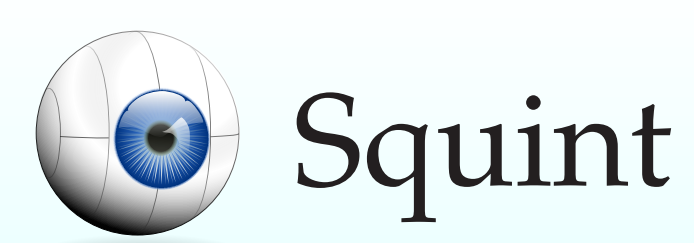
$$R_T^u := \sum_{t=1}^T (\ell_t(w_t) - \ell_t(u)), \quad \tilde{R}_T^u := \sum_{t=1}^T (w_t - u)^\top \nabla \ell_t(w_t).$$

Second-order Regret Guarantees

$$\tilde{R}_T^f \leq \sqrt{V_T^f K_T^f} \quad \text{for all } f \in \mathcal{F} \quad (1)$$

Beats worst-case regret when $V_T^{f^*} = o(T)$ and $K_T^{f^*}$ small.

Two Examples



Squint

Koolen and Van Erven [2015]

Setting Hedge Setting
 \mathcal{F} expert $k \in \{1, 2, \dots\}$

Loss Linear $w_t^\top \ell_t$

Cmplx. $K_T^k = -\ln \pi(k)$

Variance $V_T^k = \sum_{t=1}^T (w_t^\top \ell_t - \ell_t^k)^2$

Time/rd. $O(1)$ per expert



MetaGrad

Van Erven and Koolen [2016]

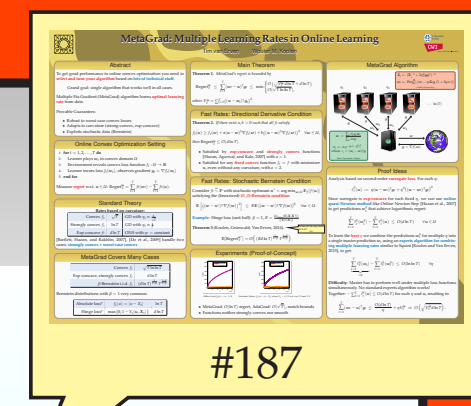
Setting Online Convex Optimization
 $u \in \mathcal{U}$

Loss Convex $\ell_t(w_t)$

Cmplx. $K_T^u = d \ln T$

Variance $V_T^u = \sum_{t=1}^T ((w_t - u)^\top \nabla \ell_t(w_t))^2$

Time/rd. $O(d^2 \ln T)$ plus projection



First Step

Consider losses $\ell \sim \mathbb{P}$ with stochastic best expert $k^* = \arg \min_k \mathbb{E}[\ell^k]$ and **gap** $\min_{k \neq k^*} \mathbb{E}[\ell^k - \ell^{k^*}] > 0$. Then second-order bound (1) implies **constant regret** $\mathbb{E}[R_T^{k^*}] = O(1)$ [Gaillard et al., 2014].

Friendly Stochastic Environments

The **Bernstein condition** [Bartlett and Mendelson, 2006] says that **variance** of excess loss is small near stochastic optimum.

Bernstein condition key to fast rates in statistical learning.

Fix $B > 0$ and $\kappa \in [0, 1]$. We say

- $\ell \sim \mathbb{P}$ are (B, κ) -Bernstein for stochastic experts if

$$\mathbb{E}[(\ell^k - \ell^{k^*})^2] \leq B \mathbb{E}[\ell^k - \ell^{k^*}]^\kappa \quad \forall k.$$

- $\ell \sim \mathbb{P}$ are (linearized) (B, κ) -Bernstein for stochastic OCO if

$$\mathbb{E} \left[((w - u^*)^\top \nabla \ell(w))^2 \right] \leq B \mathbb{E}[(w - u^*)^\top \nabla \ell(w)]^\kappa \quad \forall w.$$

See paper for extensions beyond iid.

Main Theorem

In any stochastic setting satisfying the (B, κ) -Bernstein condition, a second-order regret bound (1) implies **fast rates** both **in expectation**:

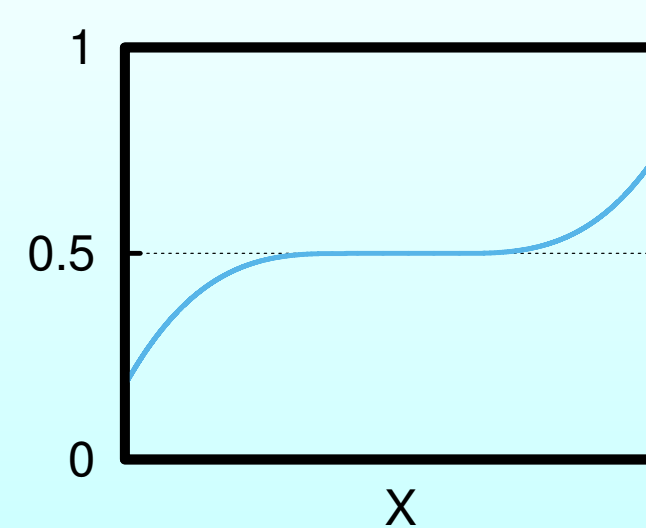
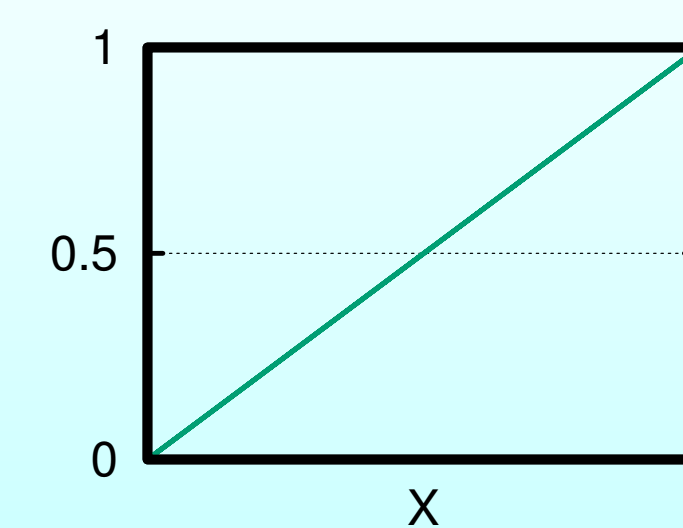
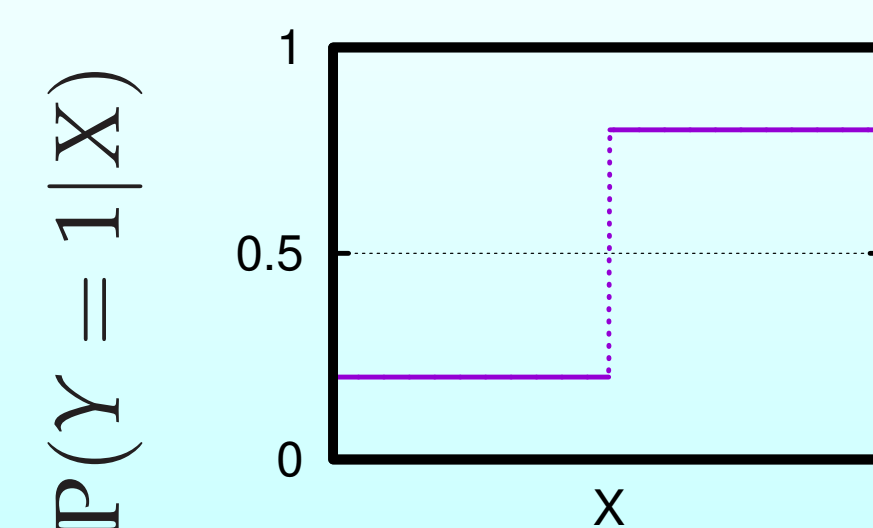
$$\mathbb{E}[R_T^{f^*}] = O \left(K_T^{\frac{1}{2-\kappa}} T^{\frac{1-\kappa}{2-\kappa}} \right),$$

and **with high probability**: for any $\delta > 0$, with probability at least $1 - \delta$,

$$R_T^{f^*} = O \left((K_T - \ln \delta)^{\frac{1}{2-\kappa}} T^{\frac{1-\kappa}{2-\kappa}} \right).$$

Inspiration: Tsybakov Margin Condition

Classification: $Y \in \{0, 1\}$. $\mathbb{P}(|\mathbb{P}(Y = 1|X) - 1/2| \leq t) \leq ct^\alpha$



Confusing case: predictors with **equal risk** but **opposite predictions**.

Hinge Loss Example

Unregularized hinge loss on unit ball.

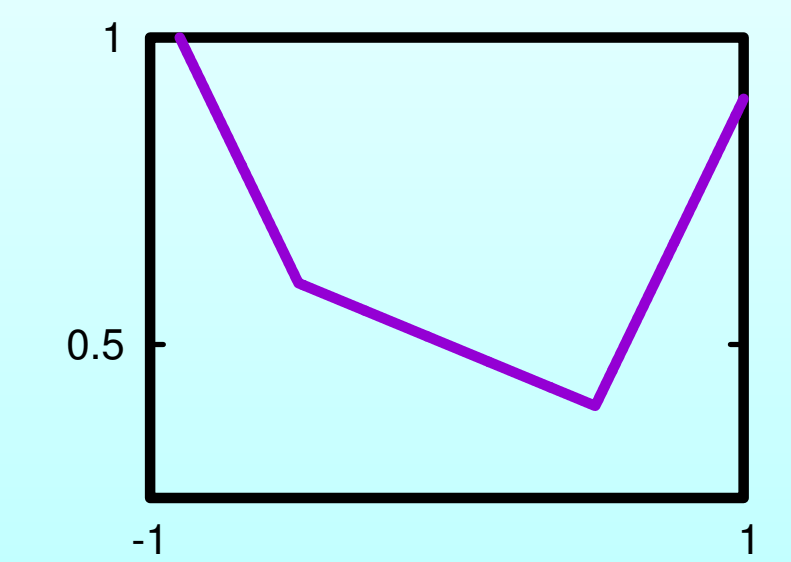
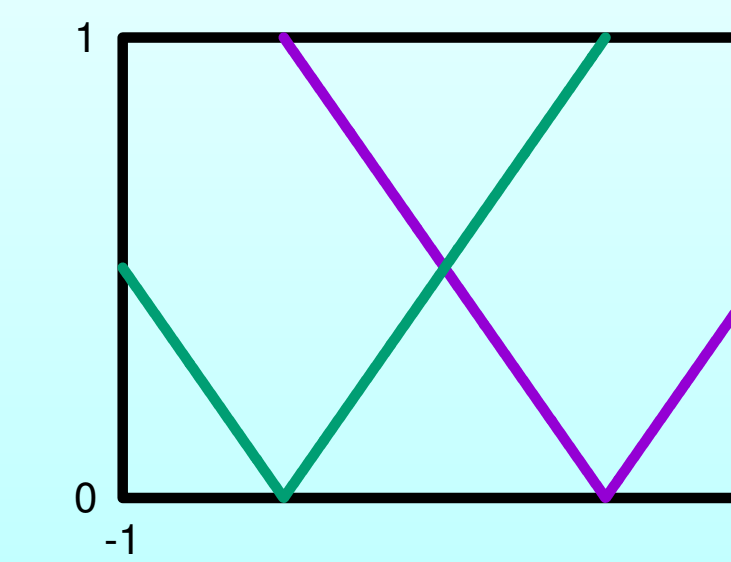
- Data $(x_t, y_t) \sim \mathbb{P}$ i.i.d.
- Hinge loss $\ell_t(u) = \max\{0, 1 - y_t x_t^\top u\}$.
- Mean $\mu = \mathbb{E}[y x]$ and second moment $D = \mathbb{E}[x x^\top]$.
- Bernstein with $\kappa = 1$ and $B = \frac{2\lambda_{\max}(D)}{\|\mu\|}$

Absolute Loss Example

Absolute loss:

$$\ell_t(u) = |u - x_t|$$

where $x_t = \pm \frac{1}{2}$ i.i.d. with probabilities $2/5$ and $3/5$.



Bernstein with $\kappa = 1$ and $B = 5$.

Proof Ideas (OCO)

In-expectation for $\kappa = 1$: Consider $\ell \sim \mathbb{P}$ with stochastic optimum $u^* = \arg \min_{u \in \mathcal{U}} \mathbb{E}[\ell(u)]$. The second-order regret bound (1) implies

$$\mathbb{E}[R_T^{u^*}] \leq \mathbb{E}[\tilde{R}_T^{u^*}] \leq \mathbb{E} \left[\sqrt{V_T^{u^*} K_T^{u^*}} \right] \leq \sqrt{\mathbb{E}[V_T^{u^*}] K_T^{u^*}}.$$

Let $x_t^u := (u - u^*)^\top \nabla \ell_t(u)$ denote the excess linearized loss of u in round t . The Bernstein condition for $\kappa = 1$ yields

$$\mathbb{E}[V_T^{u^*}] = \sum_{t=1}^T \mathbb{E}[(x_t^{u^*})^2] \leq B \sum_{t=1}^T \mathbb{E}[x_t^{u^*}] = B \mathbb{E}[\tilde{R}_T^{u^*}].$$

Combining the above two inequalities and solving for $\mathbb{E}[\tilde{R}_T^{u^*}]$ gives

$$\mathbb{E}[R_T^{u^*}] \leq B K_T^{u^*}.$$

For $\kappa < 1$: linearize ($z^\kappa = \kappa^\kappa (1 - \kappa)^{1-\kappa} \inf_{\epsilon > 0} \{e^{\kappa-1} z + \epsilon^\kappa\}$ for $z \geq 0$) to show

$$c_1 \cdot \epsilon^{1-\kappa} \mathbb{E}[V_T^{u^*}] \leq \mathbb{E}[\tilde{R}_T^{u^*}] + c_2 \cdot T \cdot \epsilon.$$

High probability: requires **sophisticated martingale argument**.