

# Boosting with the Logistic Loss is Consistent

# Boosting with the Logistic Loss is Consistent

**Boring Goal:**

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:**

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

...

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

... **why bother?**



# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

... **why bother?**

**Aspiration:**

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

... **why bother?**

**Aspiration:** Reusable techniques for similar problems.

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

... **why bother?**

**Aspiration:** Reusable techniques for similar problems.

**Strategy:**

# Boosting with the Logistic Loss is Consistent

**Boring Goal:** Statistical rates for AdaBoost with Logistic and similar strictly convex Lipschitz losses.

**Difficulty:** No regularization / constraints,  
no minimizers / strong convexity,  
linearly dependent features / singular Hessian,  
infinite dimension / nasty spectrum,  
Lipschitz  $\implies$  small Hessian on bad errors.

... **why bother?**

**Aspiration:** Reusable techniques for similar problems.

**Strategy:** Identify structure over source distribution via duality; carry it to sample.

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\inf \left\{ \text{Logistic loss of } f : f \in \text{span}(\mathcal{H}) \right\}$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\}$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\}$$
$$= \max \left\{ \text{Fermi-Dirac entropy of } p \right\}$$



## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad \left. : \right\} \end{aligned}$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\}$$

$$= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ \left. : p \text{ has capped weights,} \right.$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad \left. : p \in L^1(\mu), p \in [0, \beta] \mu\text{-a.e.}, \right\} \end{aligned}$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad \left. : p \in L^1(\mu), p \in [0, \beta] \text{ } \mu\text{-a.e.}, \right. \\ & \quad \left. p \text{ decorrelates } \mathcal{H} \text{ from } \mu \right\} \end{aligned}$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad : p \in L^1(\mu), p \in [0, \beta] \text{ } \mu\text{-a.e.}, \\ & \quad \left. \forall f \in \text{span}(\mathcal{H}) \cdot \int yf(x)p(x, y) d\mu(x, y) = 0 \right\}. \end{aligned}$$

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad \left. : p \in L^1(\mu), p \in [0, \beta] \text{ } \mu\text{-a.e.}, \right. \\ & \quad \left. \forall f \in \text{span}(\mathcal{H}) \cdot \int yf(x)p(x, y) d\mu(x, y) = 0 \right\}. \end{aligned}$$

- ▶ When optimal value positive:

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad \left. : p \in L^1(\mu), p \in [0, \beta] \text{ } \mu\text{-a.e.}, \right. \\ & \quad \left. \forall f \in \text{span}(\mathcal{H}) \cdot \int yf(x)p(x, y) d\mu(x, y) = 0 \right\}. \end{aligned}$$

- ▶ When optimal value positive:  
dual optimum certifies difficulty in every direction.

## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad \left. : p \in L^1(\mu), p \in [0, \beta] \mu\text{-a.e.}, \right. \\ & \quad \left. \forall f \in \text{span}(\mathcal{H}) \cdot \int yf(x)p(x, y) d\mu(x, y) = 0 \right\}. \end{aligned}$$

- ▶ When optimal value positive:  
dual optimum certifies difficulty in every direction.
- ▶ Difficulty in every direction  $\implies$  norm constraints.



## Nonseparable case

- ▶ When  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  is nondecreasing,  $\beta$ -Lipschitz,

$$\begin{aligned} & \inf \left\{ \int \ell(-yf(x)) d\mu(x, y) : f \in \text{span}(\mathcal{H}) \right\} \\ &= \max \left\{ - \int \ell^*(p(x, y)) d\mu(x, y) \right. \\ & \quad : p \in L^1(\mu), p \in [0, \beta] \text{ } \mu\text{-a.e.}, \\ & \quad \left. \forall f \in \text{span}(\mathcal{H}) \cdot \int yf(x)p(x, y) d\mu(x, y) = 0 \right\}. \end{aligned}$$

- ▶ When optimal value positive:  
dual optimum certifies difficulty in every direction.
- ▶ Difficulty in every direction  $\implies$  norm constraints.
- ▶ Rate  $m^{-c}$ ;  $c$  depends on  $\mathcal{H}$  and  $\mu$   $\ddot{\cdot}$ .

## Separable case

- ▶ What if optimal value zero?

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate:

$$\gamma =$$

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate:

$$\gamma = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_{\star} = 1}} \int y f(x) p(x, y) d\mu(x, y) \right.$$

:

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate:

$$\gamma = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_{\star} = 1}} \int y f(x) p(x, y) d\mu(x, y) \right.$$
$$\left. : p \in L^1(\mu), \|p\|_1 = 1, \right.$$

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate:

$$\gamma = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_{\star} = 1}} \int y f(x) p(x, y) d\mu(x, y) \right. \\ \left. : p \in L^1(\mu), \|p\|_1 = 1, p \in [0, \infty] \mu\text{-a.e.} \right\}.$$

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate:

$$\gamma = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_{\star} = 1}} \int y f(x) p(x, y) d\mu(x, y) \right. \\ \left. : p \in L^1(\mu), \|p\|_1 = 1, p \in [0, \infty] \mu\text{-a.e.} \right\}.$$

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate (adapted to Lipschitz losses):

$$\gamma_\epsilon = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_\star = 1}} \int y f(x) p(x, y) d\mu(x, y) \right. \\ \left. : p \in L^1(\mu), \|p\|_1 = 1, p \in [0, 1/\epsilon] \mu\text{-a.e.} \right\}.$$



## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate (adapted to Lipschitz losses):

$$\gamma_\epsilon = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_\star = 1}} \int y f(x) p(x, y) d\mu(x, y) \right. \\ \left. : p \in L^1(\mu), \|p\|_1 = 1, p \in [0, 1/\epsilon] \mu\text{-a.e.} \right\}.$$

- ▶ (Hi Manfred, Rocco, Satyen, Shai, ...)

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate (adapted to Lipschitz losses):

$$\gamma_\epsilon = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_\star = 1}} \int y f(x) p(x, y) d\mu(x, y) \right. \\ \left. : p \in L^1(\mu), \|p\|_1 = 1, p \in [0, 1/\epsilon] \mu\text{-a.e.} \right\}.$$

- ▶ (Hi Manfred, Rocco, Satyen, Shai, ...)
- ▶ Earlier optimal value zero  $\iff \gamma_\epsilon > 0$  for  $\epsilon > 0 \dots!$

## Separable case

- ▶ What if optimal value zero?
- ▶ Weak learning rate (adapted to Lipschitz losses):

$$\gamma_\epsilon = \inf \left\{ \sup_{\substack{f \in \text{span}(\mathcal{H}) \\ \|f\|_\star = 1}} \int y f(x) p(x, y) d\mu(x, y) \right. \\ \left. : p \in L^1(\mu), \|p\|_1 = 1, p \in [0, 1/\epsilon] \mu\text{-a.e.} \right\}.$$

- ▶ (Hi Manfred, Rocco, Satyen, Shai, ...)
- ▶ Earlier optimal value zero  $\iff \gamma_\epsilon > 0$  for  $\epsilon > 0$ ...!
- ▶  $\gamma_\epsilon$  lower bounds progress; rate  $\mathcal{O}(m^{-1/3})$ .