SVMs + Margin bounds CPSC 532S: Modern Statistical Learning Theory 2 March 2022 cs.ubc.ca/~dsuth/532S/22/

Admin

- Send in A1 regrade requests soon if I messed up
 - Either via Gradescope, or on Piazza
- A2 solutions are posted
- Later this week:
 - A3 (due in ~2 weeks after posting it)
 - Project info (with a short proposal also due in a few weeks)
 - If you don't have a group and want one, post on Piazza asap

h if I messed up azza

ng it) sal also due in a few weeks) vant one, post on Piazza asap

Linear classifiers

- Reminder: linear classifiers are
 - General case: $h(x) = \mathbb{I}(w^{\top}x + b \ge 0)$ for $x \in \mathbb{R}^d$
 - Homogenous: $h(x) = \mathbb{I}(w^{\top}x \ge 0)$

$b \ge 0$) for $x \in \mathbb{R}^d$ (0)

Linear classifiers

- Reminder: linear classifiers are
 - General case: $h(x) = \mathbb{I}(w^{\top}x + b \ge 0)$ for $x \in \mathbb{R}^d$
 - Homogenous: $h(x) = \mathbb{I}(w^{\top}x \ge 0)$
 - Can turn general case into homogenous by using $\tilde{x} = \begin{bmatrix} 1 & x \end{bmatrix} \in \mathbb{R}^{d+1}$

Linear classifiers $h(x): \frac{\chi \rightarrow \xi \sigma}{\chi \rightarrow \xi - 1/3}$ NG=sign(wtx) • General case: $h(x) = \mathbb{I}(w^{\top}x + b \ge 0)$ for $x \in \mathbb{I}(w^{\top}x + b \ge 0)$

- Reminder: linear classifiers are

 - Homogenous: $h(x) = \mathbb{I}(w^{\top}x \ge 0)$
 - Can turn general case into homogenous by using $\tilde{x} = \begin{bmatrix} 1 & x \end{bmatrix} \in \mathbb{R}^{d+1}$ • Also called *halfspace, d*efined by a hyperplane $w^{T}x + b = 0$





Linear classifiers

- Reminder: linear classifiers are
 - General case: $h(x) = \mathbb{I}(w^{\top}x + b \ge 0)$ for $x \in \mathbb{R}^d$
 - Homogenous: $h(x) = \mathbb{I}(w^{\top}x \ge 0)$
- Can turn general case into homogenous by using $\tilde{x} = \begin{bmatrix} 1 & x \end{bmatrix} \in \mathbb{R}^{d+1}$ • Also called *halfspace, d*efined by a hyperplane $w^{T}x + b = 0$ Separable case: which ERM to pick?

Margin

• Distance from x to hyperplane $\{x : w^T x + b = 0\}$ is the geometric margin

$$\rho_h(x) = \frac{|w^{\mathsf{T}}x + b|}{\|w\|_2}$$





Margin

 \bullet

$$\rho_h(x) = \frac{\|w^{\mathsf{T}}x + b\|}{\|w\|_2}$$

• Rescaling w and b doesn't change the classifier, or the margin

Distance from x to hyperplane $\{x : w^{\top}x + b = 0\}$ is the geometric margin



Max-margin classifier

• Margin for a training set is $\frac{1}{\|w\|} \min_{i \in [n]} |w^T x_i + b|$

Max-margin classifier

• Margin for a training set is $\frac{1}{\|w\|} \min_{i \in [n]} |w^T x_i + b|$ Max-margin classifier (the "hard" Support Vector Machine) is

- $\arg\max_{w,b} \frac{1}{\|w\|} \min_{i \in [n]} |w^{\mathsf{T}}x_i + b| \text{ s.t. } \forall i, \ y_i(w^{\mathsf{T}}x_i + b) > 0$

Max-margin classifier

- Generic Margin for a training set is $\frac{1}{\|w\|} \min_{i \in [n]} |w^{\top}x_i + b|$
- Max-margin classifier (the "hard" Support Vector Machine) is $\arg \max_{w,b} \frac{1}{\|w\|} \min_{i \in [n]} |w^{\top}x_i + b| \text{ s.t. } \forall i, \ y_i(w^{\top}x_i + b) > 0$
- Equivalently, as a quadratic program: $\arg\min_{w,b} \|w\|^2$ s.t. $\forall i, y_i(w^Tx_i + b) \ge 1$



• Would like to analyze 0-1 loss, but it's not Lipschitz, so we can't

- Would like to analyze 0-1 loss, but it's not Lipschitz, so we can't
 - Not just a technicality!



- Would like to analyze 0-1 loss, but it's not Lipschitz, so we can't
 - Not just a technicality!
- We'll analyze with ramp loss:

 $\mathscr{C}^{\operatorname{ramp}}(h,(x,y)) = \psi^{\operatorname{ramp}}(y)$

 $\psi^{\text{ramp}}(z) = \min\left(1, \max\left(0\right)\right)$



$$\begin{pmatrix} h(x) \\ y^{*} \\ y^{*} \\ (x,y) \end{pmatrix} = \begin{cases} 1 & \text{if } z < 0 \\ 1 - z & \text{if } 0 \le z \le 1 \\ 0 & \text{if } z > 1 \end{cases}$$

- Would like to analyze 0-1 loss, but it's not Lipschitz, so we can't
 - Not just a technicality!
- We'll analyze with ramp loss: $\mathscr{C}^{\mathrm{ramp}}(h,(x,y)) = \psi^{\mathrm{ramp}}(y)$

 $\psi^{\text{ramp}}(z) = \min\left(1, \max\left(0\right)\right)$

• ψ^{ramp} is 1-Lipschitz, bounded in [0,1]

$$h(x)) = \begin{cases} 1 & \text{if } z < 0\\ 1 - z & \text{if } 0 \le z \le 1\\ 0 & \text{if } z > 1 \end{cases}$$
], and $\ell^{0-1} \le \ell^{\text{ramp}}$; but it's *not* convex

- Would like to analyze 0-1 loss, but it's not Lipschitz, so we can't
 - Not just a technicality!
- We'll analyze with ramp loss: $\mathscr{C}^{\operatorname{ramp}}(h,(x,y)) = \psi^{\operatorname{ramp}}(y\,h(x))$

- With probability at least 1δ , it holds for all $h \in \mathcal{H}$ that

 $\psi^{\text{ramp}}(z) = \min(1, \max(0, 1 - z)) = \begin{cases} 1 & \text{if } z < 0\\ 1 - z & \text{if } 0 \le z \le 1\\ 0 & \text{if } z > 1 \end{cases}$ • ψ^{ramp} is 1-Lipschitz, bounded in [0,1], and $\ell^{0-1} \leq \ell^{\text{ramp}}$; but it's *not* convex $L_{\mathcal{D}}^{\text{ramp}}(h) \leq L_{S}^{\text{ramp}}(h) + \mathcal{P}_{\mathcal{R}}^{2} \Re_{n}\left(\mathcal{H}\right) + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$

- Would like to analyze 0-1 loss, but it's not Lipschitz, so we can't
 - Not just a technicality!
- We'll analyze with ramp loss: $\mathscr{C}^{\operatorname{ramp}}(h,(x,y)) = \psi^{\operatorname{ramp}}(y)$

 $\psi^{\text{ramp}}(z) = \min\left(1, \max\left(0\right)\right)$

- ψ^{ramp} is 1-Lipschitz, bounded in [0,1]
- With probability at least 1δ , it hold

 $L_{\mathcal{D}}^{0-1}(h) \leq L_{\mathcal{D}}^{\mathrm{ramp}}(h) \leq L_{S}^{\mathrm{ramp}}(h)$

$$h(x)) = \begin{cases} 1 & \text{if } z < 0\\ 1 - z & \text{if } 0 \le z \le 1\\ 0 & \text{if } z > 1 \end{cases}$$

], and $\ell^{0-1} \le \ell^{\text{ramp}}$; but it's *not* convex
is for all $h \in \mathcal{H}$ that
$$h^{\text{p}}(h) + \frac{2}{\checkmark} \Re_n\left(\mathcal{H}\right) + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$$

Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\operatorname{ramp}}(h) + \frac{2}{\rho} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\operatorname{ramp}}(h) + \frac{2}{\rho} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: $\Pr(y\langle w^*, x \rangle \ge 1) = 1$ $(x,y) \sim \mathcal{D}$



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\operatorname{ramp}}(h) + \frac{2}{\rho} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: Pr $(y\langle w^*, x \rangle \ge 1) = 1$ $(x,y) \sim \mathcal{D}$

- Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}}x : ||w|| \le ||w^*||\}$; assume $\mathbb{E}||x||^2 \le R^2$



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\operatorname{ramp}}(h) + \frac{2}{\rho} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: Pr $(y\langle w^*, x \rangle \ge 1) = 1$ $(x,y) \sim \mathcal{D}$

- Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}}x : \|w\| \le \|w^*\|\}$; assume $\mathbb{E}\|x\|^2 \le R^2$
 - $\mathfrak{R}_n(\mathcal{H}) \leq R \|w^*\|/\sqrt{n}$



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\operatorname{ramp}}(h) + \frac{2}{\rho} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: Pr $(y\langle w^*, x \rangle \ge 1) = 1$ $(x,y) \sim \mathcal{D}$

- Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}}x : ||w|| \le ||w^*||\}$; assume $\mathbb{E}||x||^2 \le R^2$
 - $\mathfrak{R}_n(\mathcal{H}) \leq R \|w^*\|/\sqrt{n}$
 - By assumption, $\hat{w}_S \in \mathcal{H}$ for hard SVMs



Sample complexity of <u>Hard</u> SVMs $L_{\mathcal{D}}^{0-1}(h) \le L_{S}^{\operatorname{ramp}}(h) + \frac{2}{\rho} \Re_{n}\left(\mathcal{H}\right) + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: Pr $(y\langle w^*, x \rangle \ge 1) = 1$ $(x,y) \sim \mathcal{D}$

- Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}}x : ||w|| \le ||w^*||\}$; assume $\mathbb{E}||x||^2 \le R^2$
 - $\mathfrak{R}_n(\mathcal{H}) \leq R \|w^*\|/\sqrt{n}$
 - By assumption, $\hat{w}_{S} \in \mathcal{H}$ for hard SVMs
 - Also $L_{\rm s}^{\rm ramp}(\hat{w}_{\rm S}) = 0$



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\text{ramp}}(h) + \frac{2}{\sigma} \Re_{n}\left(\mathcal{H}\right) + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: $\Pr_{(x,y) \sim \mathcal{O}}(y\langle w^*, x \rangle \geq 1) = 1$ $(x,y) \sim \mathcal{D}$ • Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}} x : \|w\| \le \|w^*\|\}$; assume $\mathbb{E}\|x\|^2 \le R^2$

- $\mathfrak{R}_n(\mathcal{H}) \leq R \|w^*\|/\sqrt{n}$
 - By assumption, $\hat{w}_{S} \in \mathscr{H}$ for hard SVMs
 - Also $L_{s}^{\text{ramp}}(\hat{w}_{S}) = 0$

$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \frac{2R\|w^{*}\|}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$$



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\text{ramp}}(h) + \frac{2}{\rho} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: $\Pr_{(x,y)\sim \mathcal{D}}(y\langle w^*, x\rangle \geq 1) = 1$ $(x,y) \sim \mathcal{D}$ • Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}} x : \|w\| \le \|w^*\|\}$; assume $\mathbb{E}\|x\|^2 \le R^2$

- $\mathfrak{R}_n(\mathcal{H}) \leq R \|w^*\|/\sqrt{n}$
 - By assumption, $\hat{w}_{S} \in \mathscr{H}$ for hard SVMs
 - Also $L_{s}^{\text{ramp}}(\hat{w}_{S}) = 0$

• Got a "slow" $1/\sqrt{n}$ rate

$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \frac{2R\|w^{*}\|}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$$

7



Sample complexity of Hard SVMs $L_{\mathcal{D}}^{0-1}(h) \leq L_{S}^{\operatorname{ramp}}(h) + \frac{2}{2} \Re_{n}(\mathcal{H}) + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$ • Assume (homogenously) realizable with margin 1: $\Pr_{(x,y)\sim \mathcal{O}}(y\langle w^*, x\rangle \geq 1) = 1$ $(x,y) \sim \mathcal{D}$ • Let $\mathscr{H} = \{x \mapsto w^{\mathsf{T}} x : ||w|| \le ||w^*||\}$; assume $\mathbb{E}||x||^2 \le R^2$

- $\mathfrak{R}_n(\mathcal{H}) \leq R \|w^*\|/\sqrt{n}$
 - By assumption, $\hat{w}_{S} \in \mathscr{H}$ for hard SVMs
 - Also $L_{s}^{\text{ramp}}(\hat{w}_{S}) = 0$

- Got a "slow" $1/\sqrt{n}$ rate
 - Can get "fast" 1/n rate with "local Rademacher complexity"

$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \frac{2R\|w^{*}\|}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}}$$



• $L_{\mathscr{D}}^{0-1}(\hat{w}_{S}) \leq \left(2R\|w^{*}\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right) / \sqrt{n}$ - but we don't know $\|w^{*}\|!$

• $L_{\mathscr{D}}^{0-1}(\hat{w}_S) \le \left(2R\|w^*\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right) / \sqrt{n}$ - but we don't know $\|w^*\|!$ • Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, $\delta_i = \frac{6\delta}{\pi^2 i^2}$ so that $\sum_{i=1}^{\infty} \delta_i = \delta$

• $L_{\mathscr{D}}^{0-1}(\hat{w}_S) \leq \left(2R\|w^*\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right)/\sqrt{n}$ - but we don't know $\|w^*\|!$

- Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, δ_i
- For each *i*, with prob at least $1 \delta_i$ hav $L_{\mathscr{D}}^{0-1}(w) \le L_{S}^{\operatorname{ramp}}(w) + \frac{2B_{i}R}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{2}{\delta_{i}}}$

$$s_i = \frac{6\delta}{\pi^2 i^2}$$
 so that $\sum_{i=1}^{\infty} \delta_i = \delta_i$
we uniformly over \mathcal{H}_i that

8

$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \left(2R\|w^{*}\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right)$$

- Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, δ_i
- For each *i*, with prob at least $1 \delta_i$ have uniformly over \mathcal{H}_i that $L_{\mathscr{D}}^{0-1}(w) \le L_{S}^{\operatorname{ramp}}(w) + \frac{2B_{i}R}{\sqrt{n}} + \sqrt{\frac{2}{2}}$ • For any w, take $i_w = \max\{1, \lceil \log_2 \frac{\|w\|}{r}\}$

 \sqrt{n} – but we don't know $||w^*||!$

$$s_i = \frac{6\delta}{\pi^2 i^2}$$
 so that $\sum_{i=1}^{\infty} \delta_i = \delta_i$

$$\frac{1}{2n}\log\frac{2}{\delta_{i}}$$

$$= \frac{1}{2n} : \text{gives } w \in \mathscr{H}_{i_{w}}$$

$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \left(2R\|w^{*}\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right)$$

- Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, δ_i
- For each *i*, with prob at least $1 \delta_i$ have uniformly over \mathcal{H}_i that $L_{\mathcal{D}}^{0-1}(w) \le L_{S}^{\operatorname{ramp}}(w) + \frac{2B_{i}R}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{2}{\delta_{i}}}$ • For any w, take $i_w = \max\{1, \lceil \log_2 \frac{\|w\|}{r} \rceil\}$

 \sqrt{n} – but we don't know $||w^*||!$

$$s = \frac{6\delta}{\pi^2 i^2}$$
 so that $\sum_{i=1}^{\infty} \delta_i = \delta_i$

(This is a corrected version of Shai+Shai's Theorem 26.14)

$$\left\{ - \right\}$$
: gives $w \in \mathcal{H}_{i_w}$



$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \left(2R\|w^{*}\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right)$$

Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, δ_i

• For each *i*, with prob at least $1 - \delta_i$ have uniformly over \mathcal{H}_i that $L_{\mathcal{D}}^{0-1}(w) \le L_{S}^{\operatorname{ramp}}(w) + \frac{2B_{i}R}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{2}{\delta_{i}}}$ • For any w, take $i_w = \max\{1, \lceil \log_2 \frac{\|w\|}{r} \rceil\}$

 $B_{i_w} \le \max\{2r, 2\|w\|\}$

 \sqrt{n} – but we don't know $||w^*||!$

$$s = \frac{6\delta}{\pi^2 i^2}$$
 so that $\sum_{i=1}^{\infty} \delta_i = \delta_i$

(This is a corrected version of Shai+Shai's Theorem 26.14)

$$\left\{ - \right\}$$
: gives $w \in \mathcal{H}_{i_w}$



$$L_{\mathcal{D}}^{0-1}(\hat{w}_S) \leq \left(2R\|w^*\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right)$$

Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, δ_i

• For each *i*, with prob at least $1 - \delta_i$ have uniformly over \mathcal{H}_i that $L_{\mathscr{D}}^{0-1}(w) \le L_{S}^{\operatorname{ramp}}(w) + \frac{2B_{i}R}{\sqrt{n}} + \sqrt{\frac{2}{2}}$ • For any w, take $i_w = \max\{1, \lceil \log_2 \frac{\|w\|}{r}\}$ $B_{i_w} \le \max\{2r, 2\|w\|\}$ $\frac{2}{s} = \frac{\pi^2}{s}$

 \sqrt{n} – but we don't know $||w^*||!$

$$s = \frac{6\delta}{\pi^2 i^2}$$
 so that $\sum_{i=1}^{\infty} \delta_i = \delta_i$

$$\frac{1}{2n} \log \frac{2}{\delta_i}$$
(This is a corrected version
Shai+Shai's Theorem 26)
$$\frac{|-|}{|}: \text{ gives } w \in \mathscr{H}_{i_w}$$

$$\left[\max\left\{1, \log_2 \frac{\|w\|}{r}\right\}\right]^2 \leq \frac{14\left(\max\left\{1, \log_2 \frac{\|w\|}{r}\right\}\right)^2}{\delta}$$



$$L_{\mathcal{D}}^{0-1}(\hat{w}_{S}) \leq \left(2R\|w^{*}\| + \sqrt{\frac{1}{2}\log\frac{1}{\delta}}\right)$$

Let $B_i = r2^i$, $\mathcal{H}_i = \{w : \|w\| \le B_i\}$, δ_i

• For each *i*, with prob at least $1 - \delta_i$ have uniformly over \mathcal{H}_i that $L_{\mathscr{D}}^{0-1}(w) \le L_{S}^{\operatorname{ramp}}(w) + \frac{2B_{i}R}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{2}{\delta_{i}}}$ • For any w, take $i_w = \max\{1, \left\lceil \log_2 \frac{\|w\|}{r} \right\rceil\}$: gives $w \in \mathcal{H}_{i_w}$ $B_{i_w} \le \max\{2r, 2\|w\|\}$ $\text{So } L_{\mathscr{D}}^{0-1}(\hat{w}_{S}) \leq \frac{4R \max\{\|\hat{w}_{S}\|, r\}}{\sqrt{n}} + \sqrt{\frac{1}{n} \max\left\{0, \log\log_{2}\frac{\|\hat{w}_{S}\|}{r}\right\}} + \frac{1}{2n}\log\frac{14}{\delta}$

 \sqrt{n} – but we don't know $||w^*||!$

$$s = \frac{6\delta}{\pi^2 i^2}$$
 so that $\sum_{i=1}^{\infty} \delta_i = \delta_i$



What if it's not linearly separable?



9



What if it's not linearly separable?

Hard SVM is

 $\underset{w,b}{\operatorname{arg\,min}} \|w\|^2 \quad s$

• Equivalent to RLM with $\mathcal{E}^{\text{hinge}}(h, (x, y)) = \max\{0, 1 - y h(x)\}$

w,b

s.t.
$$\forall i, y_i(w^{\top}x_i + b) \ge 1$$

• Soft SVM adds some slack variables ξ_i to allow violating the constraints: $\arg\min_{w,b,\xi} \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i^2 \text{ s.t. } \forall i, \ y_i(w^T x_i + b) \ge 1 - \xi_i, \ \xi_i \ge 0$ $\arg\min\,\lambda\|w\|^2 + L_{\rm c}^{\rm hinge}((w,b))$ S Yikitw


Generalization bound for Soft-SVM • Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 - y w^{\mathsf{T}} x\}$

Generalization bound for Soft-SVM • Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 - y w^{\mathsf{T}} x\}$ • This is ||x||-Lipschitz in *w*; assume $Pr(||x|| \le R) = 1$

Generalization bound for Soft-SVM • Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 - y w^{\mathsf{T}} x\}$ • This is ||x||-Lipschitz in *w*; assume $Pr(||x|| \le R) = 1$

- - Also convex, and we're using an L2 regularizer

- Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 y w^{\mathsf{T}} x\}$
 - This is ||x||-Lipschitz in *w*; assume $Pr(||x|| \le R) = 1$
 - Also convex, and we're using an L2 regularizer
- So, stability analysis from last time tells us $\mathbb{E}_{S}[L_{\infty}^{\text{hinge}}(\hat{w}_{S})] \leq L_{\infty}^{\text{hinge}}(\hat{w}_{S})$

$$\leq L_{\mathscr{D}}^{\text{hinge}}(w^*) + \lambda \|w^*\|^2 + \frac{2R^2}{\lambda n}$$

- Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 y w^{\mathsf{T}} x\}$
 - This is ||x||-Lipschitz in *w*; assume $Pr(||x|| \le R) = 1$
 - Also convex, and we're using an L2 regularizer
- So, stability analysis from last time tells us $\mathbb{E}_{S}[L_{\mathscr{D}}^{0-1}(\hat{w}_{S})] \leq \mathbb{E}_{S}[L_{\mathscr{D}}^{\text{hinge}}(\hat{w}_{S})] \leq L_{\mathscr{D}}^{\text{hinge}}(\hat{w}_{S})$

$$\leq L_{\mathscr{D}}^{\text{hinge}}(w^*) + \lambda \|w^*\|^2 + \frac{2R^2}{\lambda n}$$

- Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 y w^{\mathsf{T}} x\}$
 - This is ||x||-Lipschitz in *w*; assume $Pr(||x|| \le R) = 1$
 - Also convex, and we're using an L2 regularizer
- So, stability analysis from last time tells us $\mathbb{E}_{S}[L_{\mathscr{D}}^{0-1}(\hat{w}_{S})] \leq \mathbb{E}_{S}[L_{\mathscr{D}}^{\text{hinge}}(\hat{w}_{S})] \leq L_{\mathscr{D}}^{\text{hinge}}(\hat{w}_{S})$
 - $\mathbb{E}_{S}[L_{\mathscr{D}}^{0-1}(\hat{w}_{S})] \leq \inf_{\|w\| \leq B} L_{\mathscr{D}}^{\text{hinge}}(w)$

$$\leq L_{\mathscr{D}}^{\text{hinge}}(w^*) + \lambda \|w^*\|^2 + \frac{2R^2}{\lambda n}$$
$$) + 2RB\sqrt{\frac{2}{n}} \quad \text{if } \lambda = \frac{R}{B}\sqrt{\frac{2}{n}}$$

- Take b = 0; then $\ell^{\text{hinge}}(w, (x, y)) = \max\{0, 1 y w^{\mathsf{T}} x\}$
 - This is ||x||-Lipschitz in w; assume $Pr(||x|| \le R) = 1$
 - Also convex, and we're using an L2 regularizer
- So, stability analysis from last time tells us

$$\mathbb{E}_{S}[L_{\mathscr{D}}^{0-1}(\hat{w}_{S})] \leq \mathbb{E}_{S}[L_{\mathscr{D}}^{\text{hinge}}(\hat{w}_{S})] \leq L_{\mathscr{D}}^{\text{hinge}}(w^{*}) + \lambda \|w^{*}\|^{2} + \frac{2R^{2}}{\lambda n}$$
$$\mathbb{E}_{S}[L_{\mathscr{D}}^{0-1}(\hat{w}_{S})] \leq \inf_{\|w\| \leq B} L_{\mathscr{D}}^{\text{hinge}}(w) + 2RB\sqrt{\frac{2}{n}} \quad \text{if } \lambda = \frac{R}{B}\sqrt{\frac{2}{n}}$$

Or ramp loss analysis still works too:

$$\mathscr{L}_{\mathscr{D}}^{0-1}(\hat{w}_{S}) \leq \mathscr{L}_{\mathscr{D}}^{\mathrm{ramp}}(\hat{w}_{S}) + \frac{2RB}{\sqrt{n}} + \sqrt{\frac{1}{2n}\log\frac{1}{\delta}} \quad \text{for } \mathscr{H} = \{x \mapsto w^{\mathsf{T}}x : \|w\| \leq B$$
$$L_{\mathscr{D}}^{0-1}(\hat{w}_{S}) \leq L_{\mathscr{D}}^{\mathrm{ramp}}(\hat{w}_{S}) + \frac{4R\max\{\|\hat{w}_{S}\|, r\}}{\sqrt{n}} + \sqrt{\frac{1}{n}\max\{0, \log\log_{2}\frac{\|\hat{w}_{S}\|}{r}\}} + \frac{1}{2n}\log$$





Dimension-free rates

- Ramp loss analysis: if $\mathbb{E}||x||^2 \leq R^2$, $L^{0-1}_{\mathscr{D}}(\hat{w}) \leq L^{\mathrm{ra}}_{\mathscr{D}}$
- Soft SVM stal

$$\sum_{D}^{\operatorname{amp}}(\hat{w}) + \frac{4R \max\{\|\hat{w}_{S}\|, r\}}{\sqrt{n}} + \sqrt{\frac{1}{n}} \max\left\{0, \log \log_{2} \frac{\|\hat{w}_{S}\|}{r}\right\} + \frac{1}{2n} \log \frac{1}{2n}$$
bility analysis: if $\|x\| \le R$ a.s. and $\lambda = \frac{R}{B}\sqrt{\frac{2}{n}}$,

$$\mathbb{E}_{S}[L_{\mathcal{D}}^{0-1}(\hat{w}_{S})] \le \inf_{\|w\| \le B} L_{\mathcal{D}}^{\operatorname{hinge}}(w) + 2RB\sqrt{\frac{2}{n}}$$
d has a *d* in it! Rate only depends on $\|x\| \|w\| = \frac{\|x\|}{\operatorname{margin}}$

- Neither bound
- So we can learn in very high dimensions even infinite



Wait, but what about $VCdim(\mathcal{H})$?

- We have $VCdim(\mathcal{H}) = d$ (in the homogeneous case)
- But these analyses claim we can learn \mathcal{H} in infinite dimensions
 - ...because they assumed a margin
 - There exist (even realizable) distributions we can't learn in high-d, but they must have small margin

Why's it called a "support vector machine"?

- The ones where it's not "support" the separating hyperplane



• At convergence, hinge loss for (hopefully) most training examples will be 0





Hard SVM Duality

 $\min_{w} \frac{1}{2} \|w\|^2 \text{ s.t. } \forall i, \ y_i w^{\top} x_i \ge 1$ in class; might do next time.

Hard SVM Duality $\min_{w} \frac{1}{2} \|w\|^2 \text{ s.t. } \forall i, \ y_i w^\top x_i \ge 1 \quad = \min_{w} \max_{\alpha_i \ge 0} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i w^\top x_i)$

- called **weak duality**; \geq for **any** problem

- called **weak duality**; \geq for **any** problem but here we have strong duality:
 - it's equal

- called weak duality; \geq for **any** problem but here we have strong duality: it's equal

w optimization problem is differentiable + unconstrained



- called weak duality; \geq for **any** problem but here we have strong duality: it's equal

w optimization problem is differentiable + unconstrained setting gradient to zero:



- called **weak duality**; \geq for **any** problem but here we have strong duality: it's equal

 $W \neg$

w optimization problem is differentiable + unconstrained setting gradient to zero: n

$$+ \sum_{i=1}^{\infty} (-\alpha_i y_i x_i) = 0$$



Hard SVM Duality $\min_{w} \frac{1}{2} \|w\|^2 \text{ s.t. } \forall i, \ y_i w^\top x_i \ge 1 \quad = \min_{w} \max_{\alpha_i \ge 0} \frac{1}{2} \|w\|^2 + \sum_{i=1}^{\infty} \alpha_i (1 - y_i w^\top x_i)$ $= \max_{\alpha_i \ge 0} \min_{w} \frac{1}{2} \|w\|^2 + \sum_{i=1}^{\infty} \alpha_i (1 - y_i w^{\mathsf{T}} x_i)$ w optimization problem is differentiable + unconstrained setting gradient to zero: $w + \sum_{i=1}^{n} (-\alpha_{i} y_{i} x_{i}) = 0 \qquad w = \sum_{i=1}^{n} \alpha_{i} y_{i} x_{i}$

- called **weak duality**; \geq for **any** problem but here we have strong duality: it's equal

$$+ \sum_{i=1}^{n} (-\alpha_i y_i x_i) = 0 \qquad w = 2$$



Hard SVM Duality $\min_{w} \frac{1}{2} \|w\|^2 \text{ s.t. } \forall i, \ y_i w^\top x_i \ge 1 \quad = \min_{w} \max_{\alpha_i \ge 0} \frac{1}{2} \|w\|^2 + \sum_{i=1}^{\infty} \alpha_i (1 - y_i w^\top x_i)$ $= \max_{\alpha_i \ge 0} \min_{w} \frac{1}{2} \|w\|^2 + \sum_{i=1}^{\infty} \alpha_i (1 - y_i w^{\mathsf{T}} x_i)$ w optimization problem is differentiable + unconstrained *n* setting gradient to zero:

- called **weak duality**; \geq for **any** problem
- but here we have strong duality:
 - it's equal

 \mathcal{W} -



$$+\sum_{i=1}^{\infty} (-\alpha_i y_i x_i) = 0 \qquad w = \sum_{i=1}^{\infty} \alpha_i y_i x_i$$



- called **weak duality**; \geq for **any** problem
- but here we have strong duality:
 - it's equal



w optimization problem is differentiable + unconstrained

- setting gradient to zero: $w + \sum_{i=1}^{n} (-\alpha_{i}y_{i}x_{i}) = 0$ $w = \sum_{i=1}^{n} \alpha_{i}y_{i}x_{i}$ i=1i=1
 - $\alpha_i \geq 0$





- called **weak duality**; \geq for **any** problem
- but here we have strong duality:
 - it's equal

 \mathcal{W} -



w optimization problem is differentiable + unconstrained *n* setting gradient to zero: n

$$+\sum_{i=1}^{n} (-\alpha_i y_i x_i) = 0 \qquad w = \sum_{i=1}^{n} \alpha_i y_i x_i$$

 $\alpha_i \geq 0$

 $w = X^{\mathsf{T}} \operatorname{diag}(y) \alpha$





- called **weak duality**; \geq for **any** problem
- but here we have strong duality:
 - it's equal

 \mathcal{W} -



w optimization problem is differentiable + unconstrained setting gradient to zero: n n

$$+\sum_{i=1}^{\infty} (-\alpha_i y_i x_i) = 0 \qquad w = \sum_{i=1}^{\infty} \alpha_i y_i x_i$$

$$= \max_{\alpha_i \ge 0} \mathbf{1}^{\mathsf{T}} \alpha - \frac{1}{2} \alpha^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \alpha$$

 $w = X^{\mathsf{T}} \operatorname{diag}(y) \alpha$ $w^{\mathsf{T}} x = \alpha^{\mathsf{T}} \operatorname{diag}(y) X x$ 14



χ



- called **weak duality**; \geq for **any** problem
- but here we have strong duality:
 - it's equal

 $= \max_{\alpha_i \ge 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i y_i x_i^{\mathsf{T}} x_j y_j \alpha_j$ α_i is zero if $y_i w^{\dagger} x_i > 1$

w optimization problem is differentiable + unconstrained setting gradient to zero: n n

$$w + \sum_{i=1}^{\infty} (-\alpha_i y_i x_i) = 0 \qquad w = \sum_{i=1}^{\infty} \alpha_i y_i x_i$$

$$= \max_{\alpha_i \ge 0} \mathbf{1}^{\mathsf{T}} \alpha - \frac{1}{2} \alpha^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \alpha$$

 $w = X^{T} \operatorname{diag}(y) \alpha$ $w^{T} x = \alpha^{T} \operatorname{diag}(y) X x$ 14



χ



Soft SVM Duality $\min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 - \xi_i, \ \xi_i \ge 0$

 $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 - \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i w^\top x_i - \xi_i) - \sum_{i=1}^n \beta_i \xi_i$

- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - $= \max \min \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$ $\alpha_i, \beta_i \geq 0 \quad w, \xi$

- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - $= \max \min \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) Xw \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$ $\alpha_i, \beta_i \geq 0 \quad w, \xi$
- $2\lambda w X^{\mathsf{T}} \operatorname{diag}(y)\alpha = 0$

- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - = max min $\lambda ||w||^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$ $\alpha_i, \beta_i \geq 0 \quad w, \xi$

 $2\lambda w - X^{\mathsf{T}} \operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}} \operatorname{diag}(y)\alpha$

- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha \in \Omega} \min_{\xi} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$

 $2\lambda w - X^{\mathsf{T}} \operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}} \operatorname{diag}(y)\alpha$ $\frac{1}{n}\mathbf{1} - \alpha - \beta = 0$

- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha > 0} \min_{x} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$

 $2\lambda w - X^{\mathsf{T}}\operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}}\operatorname{diag}(y)\alpha$ $\frac{1}{n}\mathbf{1} - \alpha - \beta = 0$ $\beta = \frac{1}{n}\mathbf{1} - \alpha$



- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha_i,\beta_i \ge 0} \min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$
- $2\lambda w X^{\mathsf{T}}\operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}}\operatorname{diag}(y)\alpha$ $\frac{1}{m}\mathbf{1} \alpha \beta = 0$ $\beta = \frac{1}{m}\mathbf{1} \alpha$
 - $= \max_{\alpha_i \ge 0} \mathbf{1}^{\mathsf{T}} \alpha \frac{1}{4\lambda} \alpha^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \alpha$



- $\min_{\substack{w,\xi \\ w,\xi \ w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \ \alpha_i,\beta_i \ge 0}} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 y_i w^\top x_i \xi_i) \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha_i,\beta_i \ge 0} \min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$
- $2\lambda w X^{\mathsf{T}}\operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}}\operatorname{diag}(y)\alpha$ $\frac{1}{n}\mathbf{1} \alpha \beta = 0$ $\beta = \frac{1}{n}\mathbf{1} \alpha$
- $= \max_{\alpha_i \ge 0} \mathbf{1}^{\mathsf{T}} \alpha \frac{1}{4\lambda} \alpha^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \alpha \quad \text{ s.t. } \frac{1}{n} \ge \alpha_i$



- $\min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \sum_{i} \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{\substack{w,\xi \\ w,\xi \ \alpha_i,\beta_i \ge 0}} \max_{\substack{\lambda \|w\|^2 + \frac{1}{n} \sum_{i} \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i w^\top x_i - \xi_i)} - \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha_i,\beta_i \ge 0} \min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$
- $2\lambda w X^{\mathsf{T}}\operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}}\operatorname{diag}(y)\alpha$ $\frac{1}{n}\mathbf{1} \alpha \beta = 0$ $\beta = \frac{1}{n}\mathbf{1} \alpha$
- $= \max_{\alpha_i \ge 0} \mathbf{1}^{\mathsf{T}} \alpha \frac{1}{4\lambda} \alpha^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \alpha \quad \text{ s.t. } \frac{1}{n} \ge \alpha_i$

change variables: $2\lambda\tilde{\alpha} = \alpha$

 $= (2\lambda) \max_{0 \le \tilde{\alpha}_i \le \frac{1}{2\lambda n}} \mathbf{1}^{\mathsf{T}} \tilde{\alpha} - \frac{1}{2} \tilde{\alpha}^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \tilde{\alpha}$



- $\min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \sum_{i} \xi_i \text{ s.t. } \forall i, \ y_i w^{\mathsf{T}} x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{w,\xi} \max_{\alpha_i,\beta_i \ge 0} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i w^\top x_i - \xi_i) - \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha_i, \beta_i \ge 0} \min_{w, \xi} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) Xw \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$
- $2\lambda w X^{\mathsf{T}} \operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}} \operatorname{diag}(y)\alpha$ $\frac{1}{n}\mathbf{1} \alpha \beta = 0$ $\beta = \frac{1}{n}\mathbf{1} \alpha$
- $= \max_{\alpha_i \ge 0} \mathbf{1}^\top \alpha \frac{1}{4\lambda} \alpha^\top \operatorname{diag}(y) X X^\top \operatorname{diag}(y) \alpha \quad \text{ s.t. } \frac{1}{n} \ge \alpha_i \quad \begin{array}{c} \text{Only difference from hard} \\ \text{SVM is upper bound on } \tilde{\alpha}_i \end{array}$

change variables: $2\lambda\tilde{\alpha} = \alpha$

 $= (2\lambda) \max_{0 \le \tilde{\alpha}_i \le \frac{1}{2\lambda n}} \mathbf{1}^{\mathsf{T}} \tilde{\alpha} - \frac{1}{2} \tilde{\alpha}^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \tilde{\alpha}$



- $\min_{w,\xi} \lambda \|w\|^2 + \frac{1}{n} \sum_{i} \xi_i \quad \text{s.t. } \forall i, \ y_i w^\top x_i \ge 1 \xi_i, \ \xi_i \ge 0$ $= \min_{w,\xi} \max_{\alpha_i,\beta_i \ge 0} \lambda \|w\|^2 + \frac{1}{n} \sum_i \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i w^{\mathsf{T}} x_i - \xi_i) - \sum_{i=1}^n \beta_i \xi_i$

 - $= \max_{\alpha_i, \beta_i \ge 0} \min_{w, \xi} \lambda \|w\|^2 + \frac{1}{n} \mathbf{1}^{\mathsf{T}} \xi + \mathbf{1}^{\mathsf{T}} \alpha \alpha^{\mathsf{T}} \operatorname{diag}(y) X w \alpha^{\mathsf{T}} \xi \beta^{\mathsf{T}} \xi$
- $2\lambda w X^{\mathsf{T}}\operatorname{diag}(y)\alpha = 0$ $w = \frac{1}{2\lambda}X^{\mathsf{T}}\operatorname{diag}(y)\alpha$ $\frac{1}{n}\mathbf{1} \alpha \beta = 0$ $\beta = \frac{1}{n}\mathbf{1} \alpha$ $= \max_{\alpha_i \ge 0} \mathbf{1}^\top \alpha - \frac{1}{4\lambda} \alpha^\top \operatorname{diag}(y) X X^\top \operatorname{diag}(y) \alpha \quad \text{ s.t. } \frac{1}{n} \ge \alpha_i \quad \begin{array}{c} \text{Only difference from hard} \\ \text{SVM is upper bound on } \tilde{\alpha}_i \end{array}$ Can do with b also: change variables: $2\lambda\tilde{\alpha} = \alpha$ $= (2\lambda) \max_{0 \le \tilde{\alpha}_i \le \frac{1}{2\lambda n}} \mathbf{1}^{\mathsf{T}} \tilde{\alpha} - \frac{1}{2} \tilde{\alpha}^{\mathsf{T}} \operatorname{diag}(y) X X^{\mathsf{T}} \operatorname{diag}(y) \tilde{\alpha}$ add $\alpha^{\mathsf{T}} y = 0$ constraint, set $b = w^{\mathsf{T}} x_i - y_i$ for any SV







Karush–Kuhn–Tucker conditions

From Wikipedia, the free encyclopedia

In mathematical optimization, the Karush-Kuhn-Tucker (KKT) conditions, also known as the Kuhn-Tucker conditions, are first derivative tests (sometimes called first-order necessary conditions) for a solution in nonlinear programming to be optimal, provided that some regularity conditions are satisfied.

- Summarize the process of going through Lagrange duality for you
 - Like Lagrange multipliers, but allow inequality constraints
- Make things a lot faster once you're familiar with them
- Related conditions for when strong duality holds
 - Especially important: "Slater's condition"

FY


Practical SVM optimization

- Lot of work in the 90s/00s on efficient SVM dual optimizers

 - Also SVMlight (restrictive license)
 - These days: ThunderSVM
- Primal solvers:
 - LIBLINEAR (wrapped in scikit-learn)
 - SVMperf (restrictive license)
 - Pegasos is among the best, and you already know how to do it
 - Can still handle kernels this way, just a little less obvious

• Classic implementation: LIBSVM for the dual (wrapped in scikit-learn)

• It's just (optionally, projected) stochastic subgradient descent on hinge loss



Logistic regression and margins

- An extremely common classifier: logistic regression • $\ell^{\text{logistic}}(h, (x, y)) = \log(1 + \exp(-yh(x)))$
- On linearly separable data:

 - - (related work with other algorithms, other losses, other models...)
- On inseparable data:

 Limit of low regularization maximizes margin (Rosset, Zhu, Hastie NeurIPS-03) Gradient descent on unregularized problem does too (Soudry et al. JMLR 2018)

Gradient descent is biased towards max-margin (Ji and Telgarsky COLT-19)



- Summary Margin maximization: seems like a reasonable idea Hard SVM: exactly maximizes the margin when separable • Soft SVM: maximizes the margin with some slack equivalent to hinge loss with L2 regularization

- Two analyses:
 - Either case (actually, any linear model): ramp loss with Rademacher
 - Soft SVM: stability analysis
- Rates for both based on the margin, **not** ambient dimension
- Dual form:
 - Shows optimal w can be written as linear combination of training x_i • Can be helpful computationally if $n \ll d$ • Motivates the **kernel trick** – next time! 19