Foundations of Machine Learning Support Vector Machines

Binary Classification Problem

Training data: sample drawn i.i.d. from set $X \subseteq \mathbb{R}^N$ according to some distribution D,

 $S = ((x_1, y_1), \dots, (x_m, y_m)) \in X \times \{-1, +1\}.$

- Problem: find hypothesis $h: X \mapsto \{-1, +1\}$ in H(classifier) with small generalization error R(h).
 - choice of hypothesis set H: learning guarantees of previous lecture.

 \rightarrow linear classification (hyperplanes) if dimension N is not too large.

This Lecture

- Support Vector Machines separable case
- Support Vector Machines non-separable case
- Margin guarantees

Linear Separation



- classifiers: $H = \{ \mathbf{x} \mapsto \operatorname{sgn}(\mathbf{w} \cdot \mathbf{x} + b) : \mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R} \}.$
- geometric margin: $\rho = \min_{i \in [1,m]} \frac{|\mathbf{w} \cdot \mathbf{x}_i + b|}{\|\mathbf{w}\|}$.
- which separating hyperplane?

Optimal Hyperplane: Max. Margin

(Vapnik and Chervonenkis, 1965)



Maximum Margin



Optimization Problem

Constrained optimization:

$$\min_{\mathbf{w},b} \ \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1, i \in [1,m].$

- Properties:
 - Convex optimization.
 - Unique solution for linearly separable sample.

Optimal Hyperplane Equations

Lagrangian: for all $\mathbf{w}, b, \alpha_i \ge 0$,

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^{m} \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1].$$

KKT conditions:

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{i=1}^{m} \alpha_{i} y_{i} \mathbf{x}_{i} = 0 \iff \mathbf{w} = \sum_{i=1}^{m} \alpha_{i} y_{i} \mathbf{x}_{i}.$$

$$\nabla_{b} L = -\sum_{i=1}^{m} \alpha_{i} y_{i} = 0 \iff \sum_{i=1}^{m} \alpha_{i} y_{i} = 0.$$

$$\forall i \in [1, m], \ \alpha_{i} [y_{i} (\mathbf{w} \cdot \mathbf{x}_{i} + b) - 1] = 0.$$

Support Vectors

Complementarity conditions:

 $\alpha_i[y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1] = 0 \implies \alpha_i = 0 \lor y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1.$

Support vectors: vectors \mathbf{x}_i such that

$$\alpha_i \neq 0 \land y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1.$$

Note: support vectors are not unique.

Moving to The Dual

• Plugging in the expression of \mathbf{w} in L gives:

$$L = \frac{1}{2} \left\| \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i b + \sum_{i=1}^{m} \alpha_i \alpha_i y_i y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i b + \sum_{i=1}^{m} \alpha_i \alpha_i y_i y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i ($$



$$L = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j).$$

Equivalent Dual Opt. Problem

Constrained optimization:

$$\max_{\alpha} \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

subject to:
$$\alpha_i \ge 0 \land \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m].$$

Solution:

$$h(x) = \operatorname{sgn}\left(\sum_{i=1}^{m} \alpha_i y_i(\mathbf{x}_i \cdot \mathbf{x}) + b\right),$$

with $b = y_i - \sum_{j=1}^{m} \alpha_j y_j(\mathbf{x}_j \cdot \mathbf{x}_i)$ for any SV \mathbf{x}_i .

Leave-One-Out Error

Certain Definition: let h_S be the hypothesis output by learning algorithm L after receiving sample S of size m. Then, the leave-one-out error of L over Sis: $\widehat{=}$ 1 $\sum_{n=1}^{m}$

$$\widehat{R}_{\text{loo}}(L) = \frac{1}{m} \sum_{i=1}^{m} 1_{h_{S-\{x_i\}}(x_i) \neq f(x_i)}.$$

Property: unbiased estimate of expected error of hypothesis trained on sample of size m-1,

$$\underbrace{\mathbf{E}}_{S \sim D^{m}}[\widehat{R}_{\text{loo}}(L)] = \frac{1}{m} \sum_{i=1}^{m} \mathbf{E}[\mathbf{1}_{h_{S-\{x_{i}\}}(x_{i}) \neq f(x_{i})}] = \mathbf{E}[\mathbf{1}_{h_{S-\{x\}}(x) \neq f(x)}]$$
$$= \underbrace{\mathbf{E}}_{S' \sim D^{m-1}}[\sum_{x \sim D}[\mathbf{1}_{h_{S'}(x) \neq f(x)}]] = \underbrace{\mathbf{E}}_{S' \sim D^{m-1}}[R(h_{S'})].$$

Leave-One-Out Analysis

Theorem: let h_S be the optimal hyperplane for a sample S and let $N_{SV}(S)$ be the number of support vectors defining h_S . Then,

$$\mathop{\mathrm{E}}_{S \sim D^m} [R(h_S)] \le \mathop{\mathrm{E}}_{S \sim D^{m+1}} \left\lfloor \frac{N_{\mathrm{SV}}(S)}{m+1} \right\rfloor$$

Proof: Let $S \sim D^{m+1}$ be a sample linearly separable and let $x \in S$. If $h_{S-\{x\}}$ misclassifies x, then x must be a SV for h_S . Thus,

$$\widehat{R}_{\text{loo}}(\text{opt.-hyp.}) \le \frac{N_{\text{SV}}(S)}{m+1}$$

Notes

- Bound on expectation of error only, not the probability of error.
- Argument based on sparsity (number of support vectors). We will see later other arguments in support of the optimal hyperplanes based on the concept of margin.

This Lecture

- Support Vector Machines separable case
- Support Vector Machines non-separable case
- Margin guarantees

Support Vector Machines

(Cortes and Vapnik, 1995)

- Problem: data often not linearly separable in practice. For any hyperplane, there exists \mathbf{x}_i such that $y_i [\mathbf{w} \cdot \mathbf{x}_i + b] \geq 1.$
- dea: relax constraints using slack variables $\xi_i \ge 0$

 $y_i \left[\mathbf{w} \cdot \mathbf{x}_i + b \right] \ge 1 - \xi_i.$

Soft-Margin Hyperplanes



Support vectors: points along the margin or outliers.
 Soft margin: \(\rho = 1/||w||\).

Optimization Problem

(Cortes and Vapnik, 1995)

Constrained optimization:

$$\min_{\mathbf{w},b,\xi} \ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1 - \xi_i \land \xi_i \ge 0, i \in [1, m].$

Properties:

- $C \ge 0$ trade-off parameter.
- Convex optimization.
- Unique solution.

Notes

- Parameter C: trade-off between maximizing margin and minimizing training error. How do we determine C?
- The general problem of determining a hyperplane minimizing the error on the training set is NPcomplete (as a function of the dimension).
- Other convex functions of the slack variables could be used: this choice and a similar one with squared slack variables lead to a convenient formulation and solution.

SVM - Equivalent Problem

Optimization:

$$\min_{\mathbf{w},b} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \left(1 - y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \right)_+.$$

- Loss functions:
 - hinge loss:

$$L(h(x), y) = (1 - yh(x))_+.$$

• quadratic hinge loss:

$$L(h(x), y) = (1 - yh(x))_{+}^{2}.$$



SVMs Equations

Lagrangian: for all $\mathbf{w}, b, \alpha_i \ge 0, \beta_i \ge 0$,

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] - \sum_{i=1}^m \beta_i \xi_i.$$

• **KKT conditions:**

$$\nabla_{w}L = \mathbf{w} - \sum_{i=1}^{m} \alpha_{i}y_{i}\mathbf{x}_{i} = 0 \iff \mathbf{w} = \sum_{i=1}^{m} \alpha_{i}y_{i}\mathbf{x}_{i}.$$

$$\nabla_{b}L = -\sum_{m} \alpha_{i}y_{i} = 0 \iff \sum_{i=1}^{m} \alpha_{i}y_{i} = 0.$$

$$\nabla_{\xi_{i}}L = C - \alpha_{i} - \beta_{i} = 0 \iff \sum_{i=1}^{m} \alpha_{i}y_{i} = 0.$$

$$\forall i \in [1, m], \ \alpha_{i}[y_{i}(\mathbf{w} \cdot \mathbf{x}_{i} + b) - 1 + \xi_{i}] = 0$$

$$\beta_{i}\xi_{i} = 0.$$

Support Vectors

Complementarity conditions:

 $\alpha_i[y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0 \implies \alpha_i = 0 \lor y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 - \xi_i.$

Support vectors: vectors \mathbf{x}_i such that

$$\alpha_i \neq 0 \land y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 - \xi_i.$$

Note: support vectors are not unique.

Moving to The Dual

Plugging in the expression of w in L gives:

$$L = \frac{1}{2} \left\| \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i b + \sum_{i=1}^{m} \alpha_i \alpha_i y_i y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i b + \sum_{i=1}^{m} \alpha_i \alpha_i y_i y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i y_i y_i (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^{$$

Thus,

$$L = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j).$$

• The condition $\beta_i \ge 0$ is equivalent to $\alpha_i \le C$.

Dual Optimization Problem

Constrained optimization:

$$\max_{\alpha} \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

subject to:
$$0 \le \alpha_i \le C \land \sum_{i=1}^{\infty} \alpha_i y_i = 0, i \in [1, m].$$

Solution:

$$h(x) = \operatorname{sgn}\left(\sum_{i=1}^{m} \alpha_i y_i(\mathbf{x}_i \cdot \mathbf{x}) + b\right),$$

with $b = y_i - \sum_{j=1}^{m} \alpha_j y_j(\mathbf{x}_j \cdot \mathbf{x}_i)$ for any \mathbf{x}_i with $0 < \alpha_i < C$.

This Lecture

- Support Vector Machines separable case
- Support Vector Machines non-separable case
- Margin guarantees

High-Dimension

Learning guarantees: for hyperplanes in dimension N with probability at least $1 - \delta$,

$$R(h) \le \widehat{R}(h) + \sqrt{\frac{2(N+1)\log\frac{em}{N+1}}{m}} + \sqrt{\frac{\log\frac{1}{\delta}}{2m}}.$$

- bound is uninformative for $N \gg m$.
- but SVMs have been remarkably successful in high-dimension.
- can we provide a theoretical justification?
- analysis of underlying scoring function.

Confidence Margin

- Definition: the confidence margin of a real-valued function h at $(x, y) \in X \times Y$ is $\rho_h(x, y) = yh(x)$.
 - interpreted as the hypothesis' confidence in prediction.
 - if correctly classified coincides with |h(x)|.
 - relationship with geometric margin for linear functions $h: \mathbf{x} \mapsto \mathbf{w} \cdot \mathbf{x} + b$: for x in the sample,

 $|\rho_h(x,y)| \ge \rho_{\text{geom}} ||\mathbf{w}||.$

Confidence Margin Loss

Definition: for any confidence margin parameter $\rho > 0$ the ρ -margin loss function Φ_{ρ} is defined by



For a sample $S = (x_1, \ldots, x_m)$ and real-valued hypothesis h, the empirical margin loss is

$$\widehat{R}_{\rho}(h) = \frac{1}{m} \sum_{i=1}^{m} \Phi_{\rho}(y_i h(x_i)) \le \frac{1}{m} \sum_{i=1}^{m} 1_{y_i h(x_i) < \rho}$$

General Margin Bound

Theorem: Let H be a set of real-valued functions. Fix $\rho > 0$. For any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $h \in H$:

$$R(h) \leq \widehat{R}_{\rho}(h) + \frac{2}{\rho} \Re_{m}(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$
$$R(h) \leq \widehat{R}_{\rho}(h) + \frac{2}{\rho} \widehat{\Re}_{S}(H) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$$

Proof: Let $\widetilde{H} = \{z = (x, y) \mapsto yh(x) : h \in H\}$. Consider the family of functions taking values in [0, 1]:

$$\widetilde{\mathcal{H}} = \{ \Phi_{\rho} \circ f \colon f \in \widetilde{H} \}.$$

• By the theorem of Lecture 3, with probability at least $1-\delta$, for all $g \in \widetilde{\mathcal{H}}$,

$$\mathbf{E}[g(z)] \le \frac{1}{m} \sum_{i=1}^{m} g(z_i) + 2\mathfrak{R}_m(\widetilde{\mathcal{H}}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

Thus,

$$\operatorname{E}[\Phi_{\rho}(yh(x))] \leq \widehat{R}_{\rho}(h) + 2\mathfrak{R}_{m}(\Phi_{\rho} \circ \widetilde{H}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

- Since Φ_{ρ} is $\frac{1}{\rho}$ Lipschitz, by Talagrand's lemma, $\Re_m (\Phi_{\rho} \circ \widetilde{H}) \leq \frac{1}{\rho} \Re_m(\widetilde{H}) = \frac{1}{\rho m} \mathop{\mathrm{E}}_{\sigma,S} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i y_i h(x_i) \right] = \frac{1}{\rho} \Re_m(H).$
- Since $1_{yh(x)<0} \le \Phi_{\rho}(yh(x))$, this shows the first statement, and similarly the second one.

Rademacher Complexity of Linear Hypotheses

Theorem: Let $S \subseteq \{x : \|\mathbf{x}\| \le R\}$ be a sample of size mand let $H = \{\mathbf{x} \mapsto \mathbf{w} \cdot \mathbf{x} : \|\mathbf{w}\| \le \Lambda\}$. Then,

$$\widehat{\mathfrak{R}}_S(H) \le \sqrt{\frac{R^2 \Lambda^2}{m}}$$

Proof:

$$\begin{aligned} \widehat{\mathfrak{R}}_{S}(H) &= \frac{1}{m} \mathop{\mathbb{E}}_{\sigma} \left[\sup_{\|\mathbf{w}\| \leq \Lambda} \sum_{i=1}^{m} \sigma_{i} \mathbf{w} \cdot \mathbf{x}_{i} \right] = \frac{1}{m} \mathop{\mathbb{E}}_{\sigma} \left[\sup_{\|\mathbf{w}\| \leq \Lambda} \mathbf{w} \cdot \sum_{i=1}^{m} \sigma_{i} \mathbf{x}_{i} \right] \\ &\leq \frac{\Lambda}{m} \mathop{\mathbb{E}}_{\sigma} \left[\left\| \sum_{i=1}^{m} \sigma_{i} \mathbf{x}_{i} \right\| \right] \leq \frac{\Lambda}{m} \left[\mathop{\mathbb{E}}_{\sigma} \left[\left\| \sum_{i=1}^{m} \sigma_{i} \mathbf{x}_{i} \right\|^{2} \right] \right]^{1/2} \\ &\leq \frac{\Lambda}{m} \left[\mathop{\mathbb{E}}_{\sigma} \left[\sum_{i=1}^{m} \|\mathbf{x}_{i}\|^{2} \right] \right]^{1/2} \leq \frac{\Lambda \sqrt{mR^{2}}}{m} = \sqrt{\frac{R^{2}\Lambda^{2}}{m}}. \end{aligned}$$

Margin Bound - Linear Classifiers

Corollary: Let $\rho > 0$ and $H = {\mathbf{x} \mapsto \mathbf{w} \cdot \mathbf{x} : ||\mathbf{w}|| \le \Lambda}.$ Assume that $X \subseteq {\mathbf{x} : ||\mathbf{x}|| \le R}$. Then, for any $\delta > 0$, with probability at least $1 - \delta$, for any $h \in H$,

$$R(h) \le \widehat{R}_{\rho}(h) + 2\sqrt{\frac{R^2\Lambda^2/\rho^2}{m}} + 3\sqrt{\frac{\log\frac{2}{\delta}}{2m}}.$$

Proof: Follows directly general margin bound and bound on $\widehat{\mathfrak{R}}_{S}(H)$ for linear classifiers.

High-Dimensional Feature Space

Observations:

- generalization bound does not depend on the dimension but on the margin.
- this suggests seeking a large-margin hyperplane in a higher-dimensional feature space.
- Computational problems:
 - taking dot products in a high-dimensional feature space can be very costly.
 - solution based on kernels (next lecture).

References

- Corinna Cortes and Vladimir Vapnik, Support-Vector Networks, *Machine Learning*, 20, 1995.
- Koltchinskii, Vladimir and Panchenko, Dmitry. Empirical margin distributions and bounding the generalization error of combined classifiers. The Annals of Statistics, 30(1), 2002.
- Ledoux, M. and Talagrand, M. (1991). *Probability in Banach Spaces*. Springer, New York.
- Vladimir N.Vapnik. Estimation of Dependences Based on Empirical Data. Springer, Basederlin, 1982.
- Vladimir N.Vapnik. The Nature of Statistical Learning Theory. Springer, 1995.
- Vladimir N.Vapnik. *Statistical Learning Theory*. Wiley-Interscience, New York, 1998.

Appendix

Saddle Point

Let $(\mathbf{w}^*, b^*, \alpha^*)$ be the saddle point of the Langrangian. Multiplying both sides of the equation giving b^* by $\alpha_i^* y_i$ and taking the sum leads

$$\sum_{i=1}^m \alpha_i^* y_i b = \sum_{i=1}^m \alpha_i^* y_i^2 - \sum_{i,j=1}^m \alpha_i^* \alpha_j^* y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j).$$

• Using $y_i^2 = 1$, $\sum_{i=1}^m \alpha_i^* y_i = 0$, and $\mathbf{w}^* = \sum_{i=1}^m \alpha_i^* y_i \mathbf{x}_i$ yields

$$0 = \sum_{i=1}^{m} \alpha_i^* - \|\mathbf{w}^*\|^2.$$

Thus, the margin is also given by:

to:

$$\rho^2 = \frac{1}{\|\mathbf{w}^*\|_2^2} = \frac{1}{\|\alpha^*\|_1}.$$

Talagrand's Contraction Lemma

(Ledoux and Talagrand, 1991; pp. 112-114)

Theorem: Let $\Phi: \mathbb{R} \to \mathbb{R}$ be an *L*-Lipschitz function. Then, for any hypothesis set *H* of real-valued functions,

$$\widehat{\mathfrak{R}}_S(\Phi \circ H) \le L \,\widehat{\mathfrak{R}}_S(H).$$

Proof: fix sample $S = (x_1, \ldots, x_m)$. By definition,

$$\Re_{S}(\Phi \circ H) = \frac{1}{m} \mathop{\mathrm{E}}_{\sigma} \left[\sup_{h \in H} \sum_{i=1}^{m} \sigma_{i}(\Phi \circ h)(x_{i}) \right]$$
$$= \frac{1}{m} \mathop{\mathrm{E}}_{\sigma_{1}, \dots, \sigma_{m-1}} \left[\mathop{\mathrm{E}}_{\sigma_{m}} \left[\sup_{h \in H} u_{m-1}(h) + \sigma_{m}(\Phi \circ h)(x_{m}) \right] \right],$$
with $u_{m-1}(h) = \sum_{i=1}^{m-1} \sigma_{i}(\Phi \circ h)(x_{i}).$

Talagrand's Contraction Lemma

Now, assuming that the suprema are reached, there exist $h_1, h_2 \in H$ such that

$$\begin{split} & \underset{\sigma_m}{\mathrm{E}} \left[\sup_{h \in H} u_{m-1}(h) + \sigma_m (\Phi \circ h)(x_m) \right] \right] \\ &= \frac{1}{2} [u_{m-1}(h_1) + (\Phi \circ h_1)(x_m)] + \frac{1}{2} [u_{m-1}(h_2) - (\Phi \circ h_2)(x_m)] \\ &\leq \frac{1}{2} [u_{m-1}(h_1) + u_{m-1}(h_2) + sL(h_1(x_m) - h_2(x_m))] \\ &= \frac{1}{2} [u_{m-1}(h_1) + sLh_1(x_m)] + \frac{1}{2} [u_{m-1}(h_2) - sLh_2(x_m)] \\ &\leq \underset{\sigma_m}{\mathrm{E}} \left[\sup_{h \in H} u_{m-1}(h) + \sigma_m Lh(x_m) \right], \end{split}$$

where
$$s = sgn(h_1(x_m) - h_2(x_m))$$
.

Talagrand's Contraction Lemma

- When the suprema are not reached, the same can be shown modulo ϵ , followed by $\epsilon \rightarrow 0$.
- Proceeding similarly for other σ_i s directly leads to the result.

VC Dimension of Canonical Hyperplanes

- Theorem: Let $S \subseteq \{\mathbf{x} : \|\mathbf{x}\| \le R\}$. Then, the VC dimension d of the set of canonical hyperplanes $\{x \mapsto \operatorname{sgn}(\mathbf{w} \cdot \mathbf{x}) : \min_{x \in S} |\mathbf{w} \cdot \mathbf{x}| = 1 \land \|\mathbf{w}\| \le \Lambda\}$ verifies $d < R^2 \Lambda^2$.
- Proof: Let $\{x_1, \ldots, x_d\}$ be a set fully shattered. Then, for all $y \in \{-1, +1\}^d$, there exists w such

$$\forall i \in [1, d], 1 \leq y_i(\mathbf{w} \cdot \mathbf{x}_i).$$

Summing up the inequalities gives

$$d \leq \mathbf{w} \cdot \sum_{i=1}^{d} y_i \mathbf{x}_i \leq \|\mathbf{w}\| \|\sum_{i=1}^{d} y_i \mathbf{x}_i\| \leq \Lambda \|\sum_{i=1}^{d} y_i \mathbf{x}_i\|.$$

• Taking the expectation over $\mathbf{y} \sim U$ (uniform) yields $d \leq \Lambda \mathop{\mathrm{E}}_{\mathbf{y} \sim U} [\| \sum_{i=1}^{d} y_i \mathbf{x}_i \|] \leq \Lambda \Big[\mathop{\mathrm{E}}_{\mathbf{y} \sim U} [\| \sum_{i=1}^{d} y_i \mathbf{x}_i \|^2] \Big]^{1/2}$ (Jensen's ineq.) $= \Lambda \Big[\sum_{i,j=1}^{d} \mathrm{E}[y_i y_j] (\mathbf{x}_i \cdot \mathbf{x}_j) \Big]^{1/2}$ $= \Lambda \Big[\sum_{i=1}^{d} (\mathbf{x}_i \cdot \mathbf{x}_i) \Big]^{1/2} \leq \Lambda \Big[dR^2 \Big]^{1/2} = \Lambda R \sqrt{d}.$

• Thus, $\sqrt{d} \leq \Lambda R$.