

On-Line Learning

Nicolò Cesa-Bianchi

Università degli Studi di Milano



Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning
- 7 From mistake to risk bounds

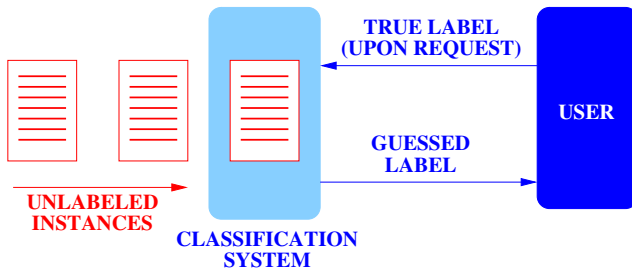


Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning
- 7 From mistake to risk bounds



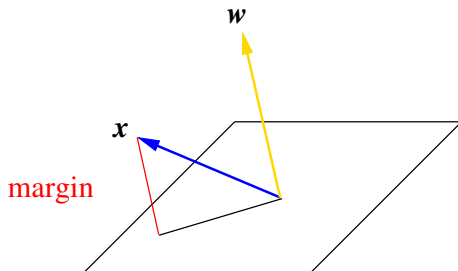
On-line classification



Linear classifiers

- Stream of data instances encoded as vectors $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^d$
- A **binary label** $y_t \in \{-1, 1\}$ associated to each \mathbf{x}_t
- A **linear classifier** $\mathbf{w}_{t-1} \in \mathbb{R}^d$ predicts label y_t of \mathbf{x}_t with

$$\hat{y}_t = \text{SGN}(\mathbf{w}_{t-1}^\top \mathbf{x}_t) \quad \mathbf{w}_{t-1} \in \mathbb{R}^d$$



On-line learning protocol

- Learning proceeds in rounds
- Initial classifier $\mathbf{w}_0 \in \mathbb{R}^d$.



On-line learning protocol

- Learning proceeds in rounds
- Initial classifier $\mathbf{w}_0 \in \mathbb{R}^d$.

On each round $t = 1, 2, \dots$

- 1 Get $\mathbf{x}_t \in \mathbb{R}^d$
- 2 Predict $\hat{y}_t = \text{SGN}(\mathbf{w}_{t-1}^\top \mathbf{x}_t)$
- 3 Observe $y_t \in \{-1, +1\}$
- 4 Update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$



On-line learning protocol

- Learning proceeds in rounds
- Initial classifier $\mathbf{w}_0 \in \mathbb{R}^d$.

On each round $t = 1, 2, \dots$

- 1 Get $\mathbf{x}_t \in \mathbb{R}^d$
- 2 Predict $\hat{y}_t = \text{SGN}(\mathbf{w}_{t-1}^\top \mathbf{x}_t)$
- 3 Observe $y_t \in \{-1, +1\}$
- 4 Update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$

- Learner observes a stream of examples (\mathbf{x}_t, y_t) $t = 1, 2, \dots$ and produces a sequence of classifiers $\mathbf{w}_0, \mathbf{w}_1, \dots$
- **Learner's performance measure:** number of classification mistakes on the stream

Remarks

- natural on certain tasks (e.g., market forecasting)
- efficient to run (scalability) and easy to code
- kernel-based learning
- strong performance guarantees (nonstochastic)
- stochastic risk bounds easy to prove (online to batch conversions)
- active learning variants when y_t comes at a cost
- tracking abilities



Summary

- 1 Linear classification
- 2 The Perceptron algorithm**
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning
- 7 From mistake to risk bounds



Perceptron algorithm

[Rosenblatt, 1952]

- Initial weight $\mathbf{w}_0 = (0, \dots, 0)$.
- Update $\mathbf{w}_t = \begin{cases} \mathbf{w}_{t-1} & \text{if } \hat{y}_t = y_t \text{ (no mistake)} \\ \mathbf{w}_{t-1} + y_t \mathbf{x}_t & \text{otherwise} \end{cases}$



Perceptron algorithm

[Rosenblatt, 1952]

- Initial weight $\mathbf{w}_0 = (0, \dots, 0)$.
- Update $\mathbf{w}_t = \begin{cases} \mathbf{w}_{t-1} & \text{if } \hat{y}_t = y_t \text{ (no mistake)} \\ \mathbf{w}_{t-1} + y_t \mathbf{x}_t & \text{otherwise} \end{cases}$

Intuition

- Want positive margin: $\hat{y}_t \neq y_t$ iff $\underbrace{y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t}_{\text{margin}} \leq 0$
- Effect of Perceptron update on margin:

$$y_t \mathbf{w}_t^\top \mathbf{x}_t = y_t (\mathbf{w}_{t-1} + y_t \mathbf{x}_t)^\top \mathbf{x}_t = y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t + \|\mathbf{x}_t\|^2$$
- Margin has increased by $\|\mathbf{x}_t\|^2$ (non-corrective update)



Linear separability

Linearly separable dataset

$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$

$$\exists \mathbf{u} \in \mathbb{R}^d \quad y_t \mathbf{u}^\top \mathbf{x}_t > 0 \quad t = 1, \dots, T$$

Note: $\|\mathbf{u}\|$ does not matter



Linear separability

Linearly separable dataset

$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$

$$\exists \mathbf{u} \in \mathbb{R}^d \quad y_t \mathbf{u}^\top \mathbf{x}_t > 0 \quad t = 1, \dots, T$$

Note: $\|\mathbf{u}\|$ does not matter

- Does the Perceptron find any such \mathbf{u} if cycled on the stream?
- What if stream is not linearly separable?



Relative loss bound

$$\sum_t m_t(\mathbf{w}_{t-1}) \leq \inf_{\mathbf{u} \in \mathbb{R}^d} \sum_t m_t(\mathbf{u}) + \text{smaller terms}$$



Relative loss bound

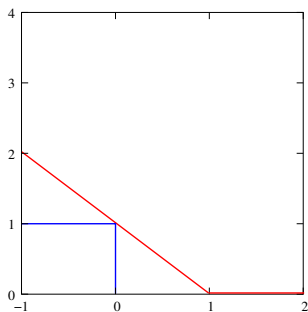
$$\sum_t m_t(\mathbf{w}_{t-1}) \leq \inf_{\mathbf{u} \in \mathbb{R}^d} \sum_t m_t(\mathbf{u}) + \text{smaller terms}$$

- Computing an hyperplane minimizing the number of misclassified examples is NP-hard



Relative loss bound

$$\sum_t m_t(\mathbf{w}_{t-1}) \leq \inf_{\mathbf{u} \in \mathbb{R}^d} \sum_t \ell_t(\mathbf{u}) + \text{smaller terms}$$

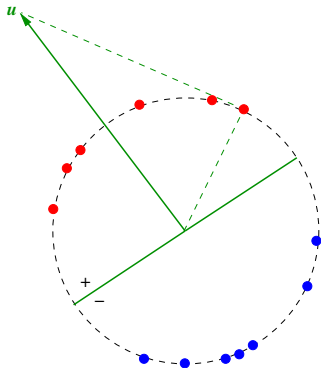


$$\underbrace{\mathbb{I}_{\{\text{SGN}(z) \neq y\}}}_{\text{mistake ind.}} \leq \underbrace{(1 - yz)_+}_{\text{hinge loss}}$$

- Computing an hyperplane minimizing the number of misclassified examples is NP-hard
- **Relax goal:** the hinge loss is a convex upper bound on the mistake indicator function
- Hinge loss grows linearly if margin is smaller than 1



Norm of the separator



- $\ell_t(\mathbf{u}) = [1 - y_t \mathbf{u}^\top \mathbf{x}_t]_+$
- Now $\|\mathbf{u}\|$ matters
- If the stream is hardly separable, then $\|\mathbf{u}\|$ has to be big to keep the hinge loss at zero



Analysis of Perceptron

- Any (nonseparable) stream, any $\mathbf{u} \in \mathbb{R}^d$, $\|\mathbf{x}_t\| = 1$
- Two proof techniques:
 1. evolution of $\|\mathbf{u} - \mathbf{w}_t\|^2$
 2. evolution of $\mathbf{u}^\top \mathbf{w}_t$



Analysis: when a mistake occurs

$$\begin{aligned} \mathbf{u}^\top \mathbf{w}_t &\geq \mathbf{u}^\top (\mathbf{w}_{t-1} + y_t \mathbf{x}_t) \\ &= \mathbf{u}^\top \mathbf{w}_{t-1} + \underbrace{y_t \mathbf{u}^\top \mathbf{x}_t}_{1 - (1 - y_t \mathbf{u}^\top \mathbf{x}_t)} \\ &\geq \mathbf{u}^\top \mathbf{w}_{t-1} + 1 - \ell_t(\mathbf{u}) \end{aligned}$$



Analysis: when a mistake occurs

$$\begin{aligned}
 \mathbf{u}^\top \mathbf{w}_t &\geq \mathbf{u}^\top (\mathbf{w}_{t-1} + y_t \mathbf{x}_t) \\
 &= \mathbf{u}^\top \mathbf{w}_{t-1} + \underbrace{y_t \mathbf{u}^\top \mathbf{x}_t}_{1 - (1 - y_t \mathbf{u}^\top \mathbf{x}_t)} \\
 &\geq \mathbf{u}^\top \mathbf{w}_{t-1} + 1 - \ell_t(\mathbf{u})
 \end{aligned}$$

Unwrapping the recurrence

$$\mathbf{u}^\top \mathbf{w}_T \geq \mathbf{u}^\top \mathbf{w}_0 + \sum_t m_t(\mathbf{w}_{t-1}) - \sum_t \ell_t(\mathbf{u})$$



When a mistake occurs (cont.)

$$\begin{aligned}\|\mathbf{w}_t\|^2 &\leq \|\mathbf{w}_{t-1} + y_t \mathbf{x}_t\|^2 \\ &= \|\mathbf{w}_{t-1}\|^2 + \|\mathbf{x}_t\|^2 + 2 \underbrace{y_t \mathbf{w}_{t-1} \mathbf{x}_t}_{\leq 0} \\ &\leq \|\mathbf{w}_{t-1}\|^2 + 1\end{aligned}$$



When a mistake occurs (cont.)

$$\begin{aligned}
 \|\mathbf{w}_t\|^2 &\leq \|\mathbf{w}_{t-1} + y_t \mathbf{x}_t\|^2 \\
 &= \|\mathbf{w}_{t-1}\|^2 + \|\mathbf{x}_t\|^2 + 2 \underbrace{y_t \mathbf{w}_{t-1} \mathbf{x}_t}_{\leq 0} \\
 &\leq \|\mathbf{w}_{t-1}\|^2 + 1
 \end{aligned}$$

Unwrapping the recurrence

$$\|\mathbf{w}_T\|^2 \leq \|\mathbf{w}_0\|^2 + \sum_t m_t(\mathbf{w}_{t-1})$$



When a mistake occurs (cont.)

$$\begin{aligned}
\|\mathbf{w}_t\|^2 &\leq \|\mathbf{w}_{t-1} + y_t \mathbf{x}_t\|^2 \\
&= \|\mathbf{w}_{t-1}\|^2 + \|\mathbf{x}_t\|^2 + 2 \underbrace{y_t \mathbf{w}_{t-1} \mathbf{x}_t}_{\leq 0} \\
&\leq \|\mathbf{w}_{t-1}\|^2 + 1
\end{aligned}$$

Unwrapping the recurrence

$$\|\mathbf{w}_T\|^2 \leq \|\mathbf{w}_0\|^2 + \sum_t m_t(\mathbf{w}_{t-1})$$

Use $\mathbf{u}^T \mathbf{w}_T \leq \|\mathbf{u}\| \|\mathbf{w}\|_T$ and solve for $\sum_t m_t(\mathbf{w}_{t-1})$



The relative mistake bound

Generalized Perceptron convergence theorem

$$\sum_t m_t(\mathbf{w}_{t-1}) \leq \inf_{\mathbf{u} \in \mathbb{R}^d} \sum_t \ell_t(\mathbf{u}) + \|\mathbf{u}\|^2 + \|\mathbf{u}\| \sqrt{\sum_t \ell_t(\mathbf{u})}$$



The relative mistake bound

Generalized Perceptron convergence theorem

$$\sum_t m_t(\mathbf{w}_{t-1}) \leq \inf_{\mathbf{u} \in \mathbb{R}^d} \sum_t \ell_t(\mathbf{u}) + \|\mathbf{u}\|^2 + \|\mathbf{u}\| \sqrt{\sum_t \ell_t(\mathbf{u})}$$

Remarks

- Oracle inequality for individual sequences
- If cycled over a linear separable $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ converges to **some** linear separator after at most $\|\mathbf{u}^*\|^2$ updates, where \mathbf{u}^* is the SVM hyperplane $\min \|\mathbf{u}\|^2 \quad \text{s.t.} \quad y_t \mathbf{u}^\top \mathbf{x}_t \geq 1 \quad \forall t$
- This implies that there always exists a linear separator expressible as a linear combination of at most $\|\mathbf{u}^*\|^2$ supports



Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams**
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning
- 7 From mistake to risk bounds



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?
- If $y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t < 1$ then update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$ so that $y_t \mathbf{w}_t^\top \mathbf{x}_t \geq 1$



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?
- If $y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t < 1$ then update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$ so that $y_t \mathbf{w}_t^\top \mathbf{x}_t \geq 1$
- Guess form of solution: $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?
- If $y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t < 1$ then update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$ so that $y_t \mathbf{w}_t^\top \mathbf{x}_t \geq 1$
- Guess form of solution: $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- $y_t (\mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t)^\top \mathbf{x}_t \geq 1$ iff $\eta_t = \frac{1}{\|\mathbf{x}_t\|^2} (1 - y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?
- If $y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t < 1$ then update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$ so that $y_t \mathbf{w}_t^\top \mathbf{x}_t \geq 1$
- Guess form of solution: $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- $y_t (\mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t)^\top \mathbf{x}_t \geq 1$ iff $\eta_t = \frac{1}{\|\mathbf{x}_t\|^2} (1 - y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$
- Equivalent to $\eta_t = \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2}$



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?
- If $y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t < 1$ then update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$ so that $y_t \mathbf{w}_t^\top \mathbf{x}_t \geq 1$
- Guess form of solution: $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- $y_t (\mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t)^\top \mathbf{x}_t \geq 1$ iff $\eta_t = \frac{1}{\|\mathbf{x}_t\|^2} (1 - y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$
- Equivalent to $\eta_t = \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2}$
- This is also the solution of the **first-order SVM**

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 \quad \text{s.t.} \quad y_t \mathbf{w}^\top \mathbf{x}_t \geq 1$$



Aggressive updates: Hildreth algorithm (1957)

- What if we want to **enforce** a large margin?
- If $y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t < 1$ then update $\mathbf{w}_{t-1} \rightarrow \mathbf{w}_t$ so that $y_t \mathbf{w}_t^\top \mathbf{x}_t \geq 1$
- Guess form of solution: $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- $y_t (\mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t)^\top \mathbf{x}_t \geq 1$ iff $\eta_t = \frac{1}{\|\mathbf{x}_t\|^2} (1 - y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$
- Equivalent to $\eta_t = \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2}$
- This is also the solution of the **first-order SVM**

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 \quad \text{s.t.} \quad y_t \mathbf{w}^\top \mathbf{x}_t \geq 1$$
- Geometrically, we project \mathbf{w}_{t-1} onto the hyperplane

$$\{\mathbf{w} \in \mathbb{R}^d : y_t \mathbf{w}^\top \mathbf{x}_t = 1\}$$



Analysis for linearly separable streams

Recall

$$\ell_t(\mathbf{w}) = [1 - \mathbf{y}_t \mathbf{w}^\top \mathbf{x}_t]_+, \text{ which implies } \ell_t(\mathbf{w}_t) = 0$$



Analysis for linearly separable streams

Recall

$$\ell_t(\mathbf{w}) = [1 - \mathbf{y}_t \mathbf{w}^\top \mathbf{x}_t]_+, \text{ which implies } \ell_t(\mathbf{w}_t) = 0$$

$$\begin{aligned} \ell_t(\mathbf{w}_{t-1}) &= \ell_t(\mathbf{w}_{t-1}) - \ell_t(\mathbf{w}_t) \\ &\leq 1 - \mathbf{y}_t \mathbf{w}_{t-1}^\top \mathbf{x}_t - 1 + \mathbf{y}_t \mathbf{w}_t^\top \mathbf{x}_t \\ &= \mathbf{y}_t (\mathbf{w}_t - \mathbf{w}_{t-1})^\top \mathbf{x}_t \\ &\leq \|\mathbf{w}_t - \mathbf{w}_{t-1}\| \|\mathbf{x}_t\| \end{aligned}$$



Analysis for linearly separable streams

Recall

$$\ell_t(\mathbf{w}) = [1 - \mathbf{y}_t \mathbf{w}^\top \mathbf{x}_t]_+, \text{ which implies } \ell_t(\mathbf{w}_t) = 0$$

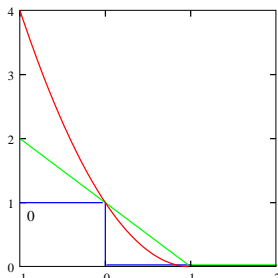
$$\begin{aligned} \ell_t(\mathbf{w}_{t-1}) &= \ell_t(\mathbf{w}_{t-1}) - \ell_t(\mathbf{w}_t) \\ &\leq 1 - \mathbf{y}_t \mathbf{w}_{t-1}^\top \mathbf{x}_t - 1 + \mathbf{y}_t \mathbf{w}_t^\top \mathbf{x}_t \\ &= \mathbf{y}_t (\mathbf{w}_t - \mathbf{w}_{t-1})^\top \mathbf{x}_t \\ &\leq \|\mathbf{w}_t - \mathbf{w}_{t-1}\| \|\mathbf{x}_t\| \end{aligned}$$

Thus we get

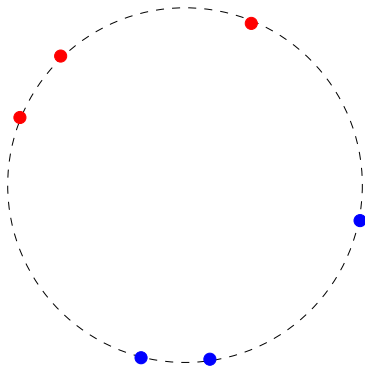
$$\ell_t(\mathbf{w}_{t-1})^2 \leq \|\mathbf{w}_t - \mathbf{w}_{t-1}\|^2 \|\mathbf{x}_t\|^2$$

Analysis (cont.)

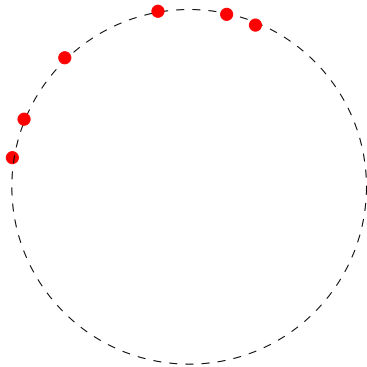
- Fix shorter separator \mathbf{u} with zero hinge loss (SVM solution)
- By the **law of cosines** $\|\mathbf{w}_t - \mathbf{w}_{t-1}\|^2 \leq \|\mathbf{w}_{t-1} - \mathbf{u}\|^2 - \|\mathbf{w}_t - \mathbf{u}\|^2$
- Thus $\sum_t m_t(\mathbf{w}_{t-1})^2 \leq \sum_t \ell_t(\mathbf{w}_{t-1})^2 \leq \|\mathbf{u}\|^2 \left(\max_t \|\mathbf{x}_t\| \right)^2$
- This is a stronger result than Perceptron convergence theorem, but it does not imply convergence to SVM solution



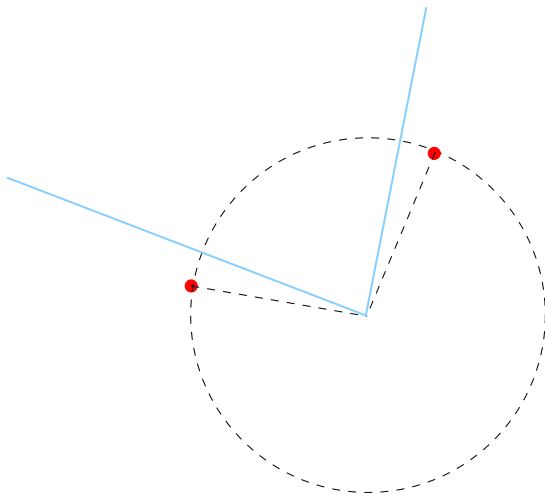
The cone of consistent hyperplanes



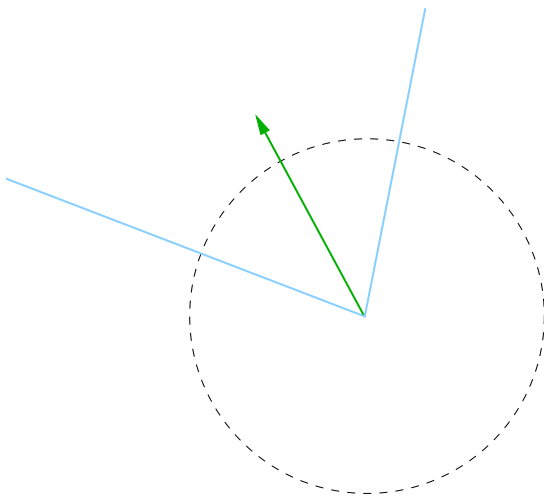
The cone of consistent hyperplanes



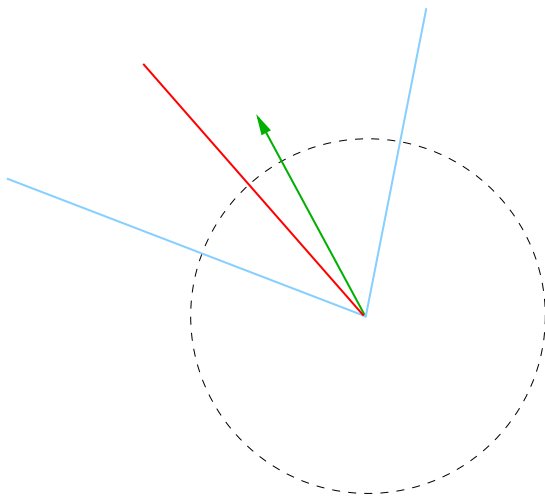
The cone of consistent hyperplanes



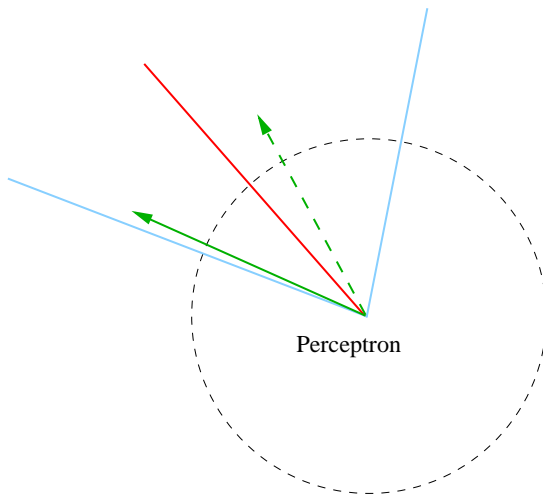
The cone of consistent hyperplanes



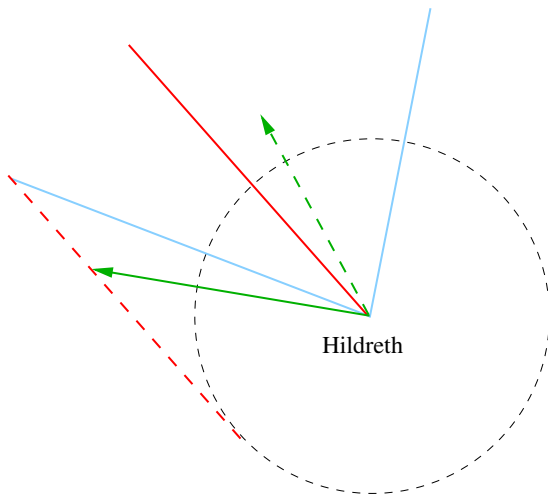
The cone of consistent hyperplanes



The cone of consistent hyperplanes



The cone of consistent hyperplanes



Mistake bounds for various updates

On any sequence of examples such that $\mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \geq 1$ with $\|\mathbf{u}\| = U$



Mistake bounds for various updates

On any sequence of examples such that $\mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \geq 1$ with $\|\mathbf{u}\| = U$

BOUND	ALGORITHM	UPDATE TIME
U^2	Perceptron	$O(d)$
U^2	Hildreth	$O(d)$
$dU \ln U$	2nd order Perceptron	$O(d^2)$
$d^2 \ln U$	Ellipsoid	$O(d^3)$
$d \ln U$	Volumetric center (U known)	$O(d^{3.5})$
$d \ln U$	Geometric center (U known)	$O(d^4)$



Mistake bounds for various updates

On any sequence of examples such that $\mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \geq 1$ with $\|\mathbf{u}\| = U$

BOUND	ALGORITHM	UPDATE TIME
U^2	Perceptron	$O(d)$
U^2	Hildreth	$O(d)$
$dU \ln U$	2nd order Perceptron	$O(d^2)$
$d^2 \ln U$	Ellipsoid	$O(d^3)$
$d \ln U$	Volumetric center (U known)	$O(d^{3.5})$
$d \ln U$	Geometric center (U known)	$O(d^4)$

Note

- U^2 optimal for fixed U
- $d \ln U$ optimal for fixed U and d



Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization**
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning
- 7 From mistake to risk bounds



Aggressive updates for nonseparable streams

Passive-aggressive algorithm

[Crammer et al. 2006]

- General update form is $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$



Aggressive updates for nonseparable streams

Passive-aggressive algorithm

[Crammer et al. 2006]

- General update form is $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- **PA-I for hinge loss** $\eta_t = \min \left\{ C, \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2} \right\}$



Aggressive updates for nonseparable streams

Passive-aggressive algorithm

[Crammer et al. 2006]

- General update form is $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- **PA-I for hinge loss** $\eta_t = \min \left\{ C, \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2} \right\}$
- This is the solution of $\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + C \ell_t(\mathbf{w})$



Aggressive updates for nonseparable streams

Passive-aggressive algorithm

[Crammer et al. 2006]

- General update form is $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- **PA-I for hinge loss** $\eta_t = \min \left\{ C, \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2} \right\}$
- This is the solution of $\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + C \ell_t(\mathbf{w})$
- **PA-II for squared hinge loss** $\eta_t = \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2 + \frac{1}{C}}$



Aggressive updates for nonseparable streams

Passive-aggressive algorithm

[Crammer et al. 2006]

- General update form is $\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_t y_t \mathbf{x}_t$
- **PA-I for hinge loss** $\eta_t = \min \left\{ C, \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2} \right\}$
- This is the solution of $\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + C \ell_t(\mathbf{w})$
- **PA-II for squared hinge loss** $\eta_t = \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2 + \frac{1}{C}}$
- This is the solution of $\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + \frac{C}{2} \ell_t(\mathbf{w})^2$



SVM and passive-aggressive

SVM primal objective

$$\text{SVM}_{\text{OPT}} = \min_{\mathbf{u}} \left(\frac{1}{2} \|\mathbf{u}\|^2 + C \sum_{t=1}^T \xi_t \right) \quad \text{s.t.} \quad y_t \mathbf{u}^\top \mathbf{x}_t \geq 1 - \xi_t \quad \forall t$$



SVM and passive-aggressive

SVM primal objective

$$\text{SVM}_{\text{OPT}} = \min_{\mathbf{u}} \left(\frac{1}{2} \|\mathbf{u}\|^2 + C \sum_{t=1}^T \xi_t \right) \quad \text{s.t.} \quad y_t \mathbf{u}^\top \mathbf{x}_t \geq 1 - \xi_t \quad \forall t$$

SVM dual Lagrange function

$$D_{\text{svm}}(\alpha_1, \dots, \alpha_T) = -\frac{1}{2} \left\| \sum_t \alpha_t \mathbf{x}_t \right\|^2 + \sum_t \alpha_t \quad 0 \leq \alpha_t \leq C \quad \forall t$$



SVM and passive-aggressive

SVM primal objective

$$\text{SVM}_{\text{OPT}} = \min_{\mathbf{u}} \left(\frac{1}{2} \|\mathbf{u}\|^2 + C \sum_{t=1}^T \xi_t \right) \quad \text{s.t.} \quad y_t \mathbf{u}^\top \mathbf{x}_t \geq 1 - \xi_t \quad \forall t$$

SVM dual Lagrange function

$$D_{\text{svm}}(\alpha_1, \dots, \alpha_T) = -\frac{1}{2} \left\| \sum_t \alpha_t \mathbf{x}_t \right\|^2 + \sum_t \alpha_t \quad 0 \leq \alpha_t \leq C \quad \forall t$$

Duality

For all $\alpha'_1, \dots, \alpha'_T \geq 0$

$$D_{\text{svm}}(\alpha'_1, \dots, \alpha'_T) \leq \sup_{\alpha} D_{\text{svm}}(\alpha) = \text{SVM}_{\text{OPT}}$$

SVM and passive-aggressive (cont.)

Following [Shalev-Shwartz and Singer, 2006]

$$\begin{aligned}D_{\text{svm}}(0, \dots, 0) &= 0 \\D_{\text{svm}}(\alpha_1, \dots, \alpha_T) &= \sum_{t=1}^T \left(D_{\text{svm}}(\alpha_1, \dots, \alpha_t, 0, \dots, 0) \right. \\&\quad \left. - D_{\text{svm}}(\alpha_1, \dots, \alpha_{t-1}, 0, \dots, 0) \right)\end{aligned}$$



SVM and passive-aggressive (cont.)

Following [\[Shalev-Shwartz and Singer, 2006\]](#)

$$D_{\text{svm}}(0, \dots, 0) = 0$$

$$D_{\text{svm}}(\alpha_1, \dots, \alpha_T) = \sum_{t=1}^T \left(D_{\text{svm}}(\alpha_1, \dots, \alpha_t, 0, \dots, 0) - D_{\text{svm}}(\alpha_1, \dots, \alpha_{t-1}, 0, \dots, 0) \right)$$

$$D_{\text{svm}}(\alpha_1, \dots, \alpha_t, 0, \dots, 0) - D_{\text{svm}}(\alpha_1, \dots, \alpha_{t-1}, 0, \dots, 0)$$

$$= -\frac{1}{2} \alpha_t^2 \|\mathbf{x}_t\|^2 + \alpha_t (1 - \mathbf{y}_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$$

where $\mathbf{w}_{t-1} = \sum_{s=1}^{t-1} \alpha_s \mathbf{y}_s \mathbf{x}_s$

SVM and passive-aggressive (cont.)

PA-I primal objective

$$\min_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + C\xi_t \right) \quad \text{s.t.} \quad \mathbf{y}_t \mathbf{w}^\top \mathbf{x}_t \geq 1 - \xi_t$$



SVM and passive-aggressive (cont.)

PA-I primal objective

$$\min_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + C\xi_t \right) \quad \text{s.t.} \quad \mathbf{y}_t \mathbf{w}^\top \mathbf{x}_t \geq 1 - \xi_t$$

PA-I dual Lagrange function

$$D_t(\alpha) = -\frac{1}{2} \alpha^2 \|\mathbf{x}_t\|^2 + \alpha(1 - \mathbf{y}_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$$



SVM and passive-aggressive (cont.)

PA-I primal objective

$$\min_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-1}\|^2 + C\xi_t \right) \quad \text{s.t.} \quad \mathbf{y}_t \mathbf{w}^\top \mathbf{x}_t \geq 1 - \xi_t$$

PA-I dual Lagrange function

$$D_t(\alpha) = -\frac{1}{2} \alpha^2 \|\mathbf{x}_t\|^2 + \alpha(1 - \mathbf{y}_t \mathbf{w}_{t-1}^\top \mathbf{x}_t)$$

PA-I learning rate is maximizer of dual objective

$$\eta_t = \operatorname{argmin} \left\{ C, \frac{\ell_t(\mathbf{w}_{t-1})}{\|\mathbf{x}_t\|^2} \right\} = \operatorname{argmax}_{0 \leq \alpha \leq C} D_t(\alpha) \stackrel{\text{def}}{=} \alpha_t$$



Mistake bounds for PA-I

We note that

$$D_t(\alpha_t) = D_{\text{svm}}(\alpha_1, \dots, \alpha_t, 0, \dots, 0) - D_{\text{svm}}(\alpha_1, \dots, \alpha_{t-1}, 0, \dots, 0)$$



Mistake bounds for PA-I

We note that

$$D_t(\alpha_t) = D_{\text{svm}}(\alpha_1, \dots, \alpha_t, 0, \dots, 0) - D_{\text{svm}}(\alpha_1, \dots, \alpha_{t-1}, 0, \dots, 0)$$

This shows

$$\sum_{t=1}^T D_t(\alpha_t) \leq \text{SVM}_{\text{OPT}} = \min_{\mathbf{u}} \left(\frac{1}{2} \|\mathbf{u}\|^2 + C \sum_{t=1}^T \ell_t(\mathbf{u}) \right)$$



Mistake bounds for PA-I

We note that

$$D_t(\alpha_t) = D_{\text{svm}}(\alpha_1, \dots, \alpha_t, 0, \dots, 0) - D_{\text{svm}}(\alpha_1, \dots, \alpha_{t-1}, 0, \dots, 0)$$

This shows

$$\sum_{t=1}^T D_t(\alpha_t) \leq \text{SVM}_{\text{OPT}} = \min_{\mathbf{u}} \left(\frac{1}{2} \|\mathbf{u}\|^2 + C \sum_{t=1}^T \ell_t(\mathbf{u}) \right)$$

- We can use this inequality to prove mistake bounds for PA-I
- **Proof idea:** lower bound $D_t(\alpha_t)$ at mistakes



Proof of mistake bound for PA-I

Assume for simplicity $\|\mathbf{x}_t\| = 1$



Proof of mistake bound for PA-I

Assume for simplicity $\|\mathbf{x}_t\| = 1$

Lemma

$$y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t \leq 0 \quad \text{implies} \quad D_t(\alpha_t) \geq \min\{C, \ell_t(\mathbf{w}_{t-1})\} \frac{\ell_t(\mathbf{w}_{t-1})}{2}$$



Proof of mistake bound for PA-I

Assume for simplicity $\|\mathbf{x}_t\| = 1$

Lemma

$$y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t \leq 0 \quad \text{implies} \quad D_t(\alpha_t) \geq \min\{C, \ell_t(\mathbf{w}_{t-1})\} \frac{\ell_t(\mathbf{w}_{t-1})}{2}$$

This implies

$$\min\{C, 1\} m_t(\mathbf{w}_{t-1}) \leq \min\{C, \underbrace{\ell_t(\mathbf{w}_{t-1})}_{\geq 1}\} \underbrace{\ell_t(\mathbf{w}_{t-1})}_{\geq m_t(\mathbf{w}_{t-1})} \leq 2 D_t(\alpha_t)$$



Proof of mistake bound for PA-I

Assume for simplicity $\|\mathbf{x}_t\| = 1$

Lemma

$$y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t \leq 0 \quad \text{implies} \quad D_t(\alpha_t) \geq \min\{C, \ell_t(\mathbf{w}_{t-1})\} \frac{\ell_t(\mathbf{w}_{t-1})}{2}$$

This implies

$$\min\{C, 1\} m_t(\mathbf{w}_{t-1}) \leq \min\{C, \underbrace{\ell_t(\mathbf{w}_{t-1})}_{\geq 1}\} \underbrace{\ell_t(\mathbf{w}_{t-1})}_{\geq m_t(\mathbf{w}_{t-1})} \leq 2 D_t(\alpha_t)$$

Theorem

[Crammer et al., 2006]

$$\sum_{t=1}^T m_t(\mathbf{w}_{t-1}) \leq \max\left\{\frac{1}{C}, 1\right\} \inf_{\mathbf{u} \in \mathbb{R}^d} \left(\|\mathbf{u}\|^2 + 2C \sum_{t=1}^T \ell_t(\mathbf{u}) \right)$$

Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning**
- 6 Online SVM and active learning
- 7 From mistake to risk bounds



On-line learning with kernels

- Feature map $\phi : \mathbb{R}^d \rightarrow \text{RKHS}$

On-line learning with kernels

- Feature map $\phi : \mathbb{R}^d \rightarrow \text{RKHS}$
- Kernel $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$

On-line learning with kernels

- Feature map $\phi : \mathbb{R}^d \rightarrow \text{RKHS}$
- Kernel $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$
- Assume a linear algorithm learns \mathbf{w} such that

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_{t_i}$$

On-line learning with kernels

- Feature map $\phi : \mathbb{R}^d \rightarrow \text{RKHS}$
- Kernel $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$
- Assume a linear algorithm learns \mathbf{w} such that

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_{t_i}$$

- Then we can learn $\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_{t_i})$ in the RKHS because

$$\begin{aligned} \text{SGN}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle) &= \text{SGN} \left(\sum_i y_{t_i} \langle \phi(\mathbf{x}_{t_i}), \phi(\mathbf{x}) \rangle \right) \\ &= \text{SGN} \left(\sum_i y_{t_i} K(\mathbf{x}_{t_i}, \mathbf{x}) \right) \end{aligned}$$

Kernel Perceptron

Start with empty cache \mathcal{L} of examples

Loop:



Kernel Perceptron

Start with empty cache \mathcal{L} of examples

Loop:

- 1 Read next instance x_t



Kernel Perceptron

Start with empty cache \mathcal{L} of examples

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} K(\mathbf{v}, \mathbf{x}_t)\right)$



Kernel Perceptron

Start with empty cache \mathcal{L} of examples

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} K(\mathbf{v}, \mathbf{x}_t)\right)$
- 3 Obtain true label y_t



Kernel Perceptron

Start with empty cache \mathcal{L} of examples

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} K(\mathbf{v}, \mathbf{x}_t)\right)$
- 3 Obtain true label y_t
- 4 If $\hat{y}_t \neq y_t$ (**mistake**) then store new **support** $(y_t \mathbf{x}_t)$ in \mathcal{L}



Kernel Perceptron

Start with empty cache \mathcal{L} of examples

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} K(\mathbf{v}, \mathbf{x}_t)\right)$
- 3 Obtain true label y_t
- 4 If $\hat{y}_t \neq y_t$ (**mistake**) then store new **support** $(y_t \mathbf{x}_t)$ in \mathcal{L}

Mistake bounds hold in the whole **RKHS**



Memory bounded learning

Can we control the rate of mistakes when at most $B < \infty$ supports are used?



Memory bounded learning

Can we control the rate of mistakes when at most $B < \infty$ supports are used?

Fact

Using at most B supports, any learner makes an unbounded number of mistakes on a sequence that is perfectly classified by some $\mathbf{u} \in \mathbb{R}^d$ with zero hinge loss and $\|\mathbf{u}\| = \sqrt{B+1}$



Memory bounded learning

Can we control the rate of mistakes when at most $B < \infty$ supports are used?

Fact

Using at most B supports, any learner makes an unbounded number of mistakes on a sequence that is perfectly classified by some $\mathbf{u} \in \mathbb{R}^d$ with zero hinge loss and $\|\mathbf{u}\| = \sqrt{B+1}$

- Thus $B \geq U^2$ is necessary to compete against \mathbf{u} of length U



Memory bounded learning

Can we control the rate of mistakes when at most $B < \infty$ supports are used?

Fact

Using at most B supports, any learner makes an unbounded number of mistakes on a sequence that is perfectly classified by some $\mathbf{u} \in \mathbb{R}^d$ with zero hinge loss and $\|\mathbf{u}\| = \sqrt{B+1}$

- Thus $B \geq U^2$ is necessary to compete against \mathbf{u} of length U
- Can we compete against any \mathbf{u} with $\|\mathbf{u}\| \leq U$ using $B = (1 + \epsilon)U^2$ supports?



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

Loop:

- 1 Read next instance x_t



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} \mathbf{v}^\top \mathbf{x}_t\right)$



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} \mathbf{v}^\top \mathbf{x}_t\right)$
- 3 Obtain true label \mathbf{y}_t



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{v \in \mathcal{L}} \mathbf{v}^\top \mathbf{x}_t\right)$
- 3 Obtain true label y_t
- 4 If $\hat{y}_t \neq y_t$ then:



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}\left(\sum_{v \in \mathcal{L}} \mathbf{v}^\top \mathbf{x}_t\right)$
- 3 Obtain true label y_t
- 4 If $\hat{y}_t \neq y_t$ then:
 - 1 If $|\mathcal{L}| = B$, then **throw away a random support** from \mathcal{L}



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} \mathbf{v}^\top \mathbf{x}_t\right)$
- 3 Obtain true label \mathbf{y}_t
- 4 If $\hat{\mathbf{y}}_t \neq \mathbf{y}_t$ then:
 - 1 If $|\mathcal{L}| = B$, then **throw away a random support** from \mathcal{L}
 - 2 Add $\mathbf{y}_t \mathbf{x}_t$ to \mathcal{L}



A randomized perceptron

Randomized Budget Perceptron

Parameter: size B of cache for supports

Start with empty cache \mathcal{L}

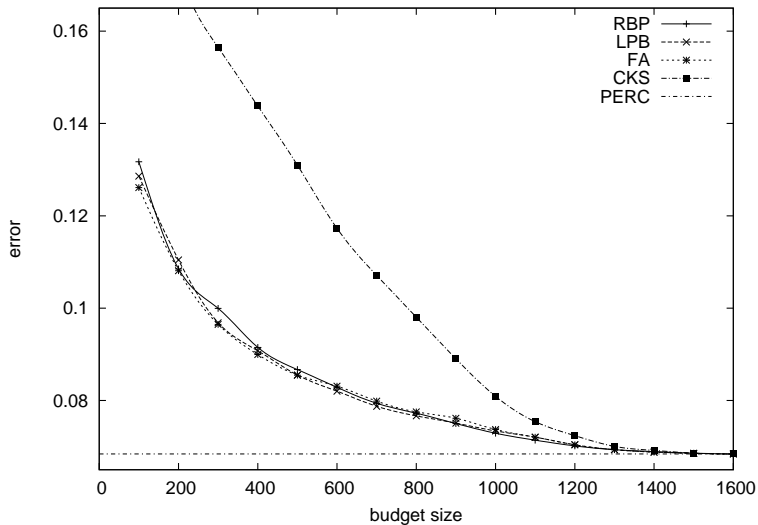
Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}\left(\sum_{\mathbf{v} \in \mathcal{L}} \mathbf{v}^\top \mathbf{x}_t\right)$
- 3 Obtain true label \mathbf{y}_t
- 4 If $\hat{\mathbf{y}}_t \neq \mathbf{y}_t$ then:
 - 1 If $|\mathcal{L}| = B$, then **throw away a random support** from \mathcal{L}
 - 2 Add $\mathbf{y}_t \mathbf{x}_t$ to \mathcal{L}

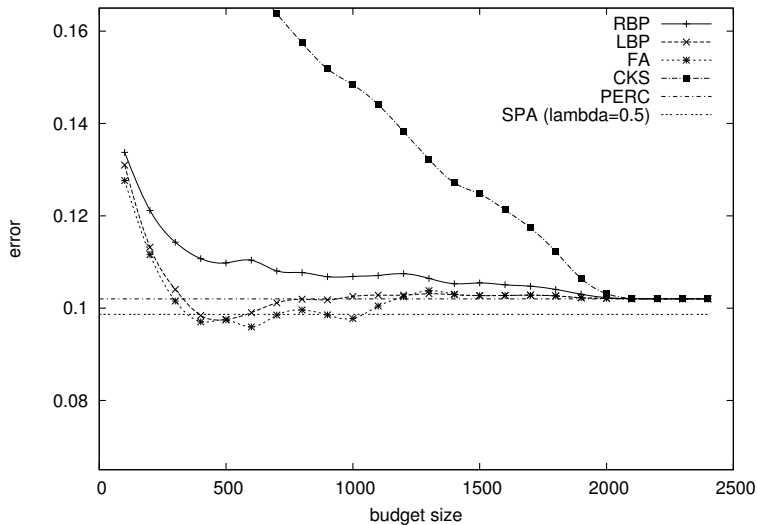
Result:

Bound on mistakes scales roughly with $1 + 1/\epsilon$
 (when using $B = (1 + \epsilon)U^2$ supports)

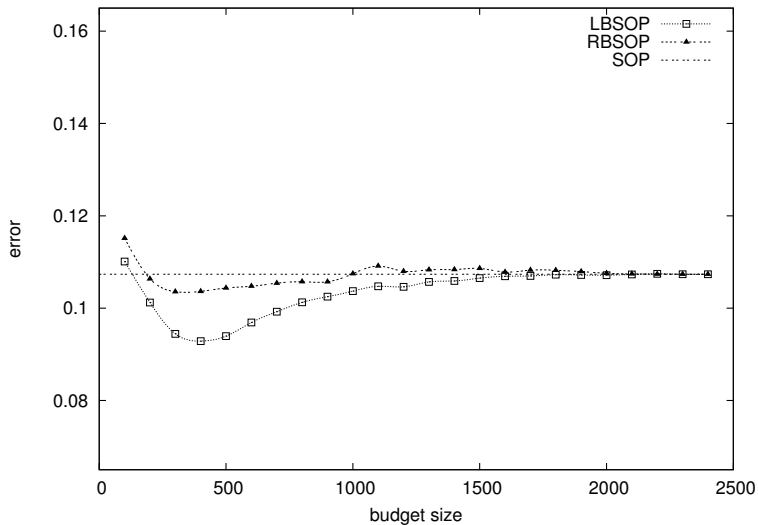
Empirical performance — stationary



Empirical performance — nonstationary



Empirical performance 2nd order — nonstationary



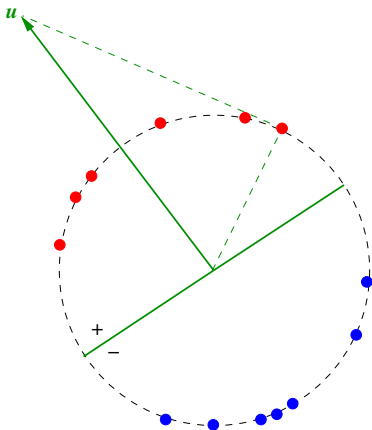
Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning**
- 7 From mistake to risk bounds



Online approximation of SVM hyperplane

The SVM hyperplane is the shortest \mathbf{u} such that $\mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \geq 1$ for all t



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

Loop:

- 1 Read next instance x_t



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 3 Obtain true label \mathbf{y}_t



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 3 Obtain true label \mathbf{y}_t
- 4 If margin smaller than $c(1 - \alpha)/\sqrt{k}$ then:



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 3 Obtain true label y_t
- 4 If margin smaller than $c(1 - \alpha)/\sqrt{k}$ then:
 - 1 $\mathbf{w}' = \mathbf{w} + y_t \mathbf{x}_t / \sqrt{k}$



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict y_t with $\hat{y}_t = \text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 3 Obtain true label y_t
- 4 If margin smaller than $c(1 - \alpha)/\sqrt{k}$ then:
 - 1 $\mathbf{w}' = \mathbf{w} + y_t \mathbf{x}_t / \sqrt{k}$
 - 2 $\mathbf{w} = \mathbf{w}' / \|\mathbf{w}'\|$ $k \leftarrow k + 1$



Online approximation of SVM hyperplane (cont.)

The ALMA algorithm

Parameter: $0 < \alpha \leq 1$

Set mistake counter $k = 1$

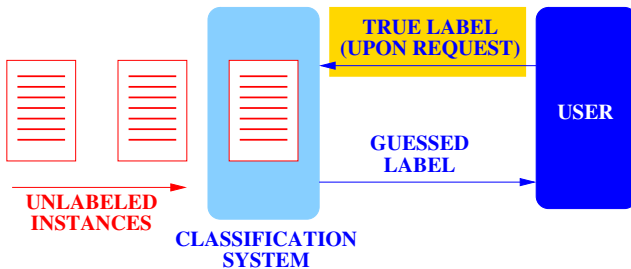
Loop:

- 1 Read next instance \mathbf{x}_t
- 2 Predict \mathbf{y}_t with $\hat{\mathbf{y}}_t = \text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 3 Obtain true label \mathbf{y}_t
- 4 If margin smaller than $c(1 - \alpha)/\sqrt{k}$ then:
 - 1 $\mathbf{w}' = \mathbf{w} + \mathbf{y}_t \mathbf{x}_t / \sqrt{k}$
 - 2 $\mathbf{w} = \mathbf{w}' / \|\mathbf{w}'\| \quad k \leftarrow k + 1$

Result

Finds separating \mathbf{u} with $\|\mathbf{u}\| \leq \|\mathbf{u}_{\text{SVM}}\| / (1 - \alpha)$ after at most $(\|\mathbf{u}_{\text{SVM}}\| / \alpha)^2$ updates

Selective sampling



A selective sampling classifier

- 1 Classify next instance \mathbf{x}_t with $\text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$



A selective sampling classifier

- 1 Classify next instance \mathbf{x}_t with $\text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 2 If $|\mathbf{w}^\top \mathbf{x}_t| \leq \|\mathbf{x}_t\| \sqrt{\frac{c \ln t}{N_t}}$ then query label y_t of \mathbf{x}_t

N_t = number of labels sampled so far



A selective sampling classifier

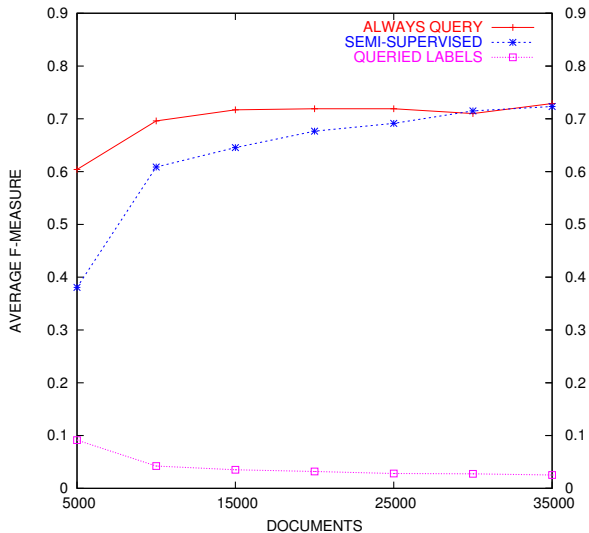
- 1 Classify next instance \mathbf{x}_t with $\text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- 2 If $|\mathbf{w}^\top \mathbf{x}_t| \leq \|\mathbf{x}_t\| \sqrt{\frac{c \ln t}{N_t}}$ then query label y_t of \mathbf{x}_t
- 3 If label queried then use (\mathbf{x}_t, y_t) to update \mathbf{w}

N_t = number of labels sampled so far

\mathbf{w} updated with the **2nd order Perceptron update rule**



Empirical performance on RCV1



Summary

- 1 Linear classification
- 2 The Perceptron algorithm
- 3 Mistake bounds for separable streams
- 4 Online learning and convex optimization
- 5 Kernel-based on-line learning
- 6 Online SVM and active learning
- 7 From mistake to risk bounds**



Statistical learning theory

- Linear classifiers $H(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x})$



Statistical learning theory

- Linear classifiers $H(\mathbf{x}) = \text{SGN}(\mathbf{w}^\top \mathbf{x})$
- Examples (\mathbf{x}_t, y_t) are i.i.d. according to a fixed and unknown probability distribution on $\mathbb{R}^d \times \{-1, +1\}$



Statistical learning theory

- Linear classifiers $H(\mathbf{x}) = \text{SGN}(\mathbf{w}^\top \mathbf{x})$
- Examples (\mathbf{x}_t, y_t) are i.i.d. according to a fixed and unknown probability distribution on $\mathbb{R}^d \times \{-1, +1\}$
- $\text{risk}(H) = \mathbb{P}(H(\mathbf{x}) \neq y)$



Statistical learning theory

- Linear classifiers $H(\mathbf{x}) = \text{SGN}(\mathbf{w}^\top \mathbf{x})$
- Examples (\mathbf{x}_t, y_t) are i.i.d. according to a fixed and unknown probability distribution on $\mathbb{R}^d \times \{-1, +1\}$
- $\text{risk}(H) = \mathbb{P}(H(\mathbf{x}) \neq y)$
- Learning algorithm

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \longrightarrow \boxed{A} \longrightarrow \hat{H} : \mathbb{R}^d \rightarrow \{-1, +1\}$$

\hat{H} is (random) hypothesis output by learner



The ensemble of hypotheses

- Run an incremental learner on the training set



The ensemble of hypotheses

- Run an incremental learner on the training set
- Everytime $H(\mathbf{x}_t) \neq y_t$, H is changed by the update rule



The ensemble of hypotheses

- Run an incremental learner on the training set
- Everytime $H(x_t) \neq y_t$, H is changed by the update rule
- This process generates an **ensemble of classifiers**

$$H_0, H_1, \dots, H_n$$



The ensemble of hypotheses

- Run an incremental learner on the training set
- Everytime $H(x_t) \neq y_t$, H is changed by the update rule
- This process generates an **ensemble of classifiers**

$$H_0, H_1, \dots, H_n$$

Goals

- 1 Bound the **average risk of the ensemble** in terms of the size of the ensemble



The ensemble of hypotheses

- Run an incremental learner on the training set
- Everytime $H(x_t) \neq y_t$, H is changed by the update rule
- This process generates an **ensemble of classifiers**

$$H_0, H_1, \dots, H_n$$

Goals

- 1 Bound the **average risk of the ensemble** in terms of the size of the ensemble
- 2 Find an element of the ensemble whose risk is close to the ensemble average



Step 1: bound the average risk

The difference

$$\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq y_t\}}$$

is a **martingale difference sequence** because

$$\mathbb{E} \left[\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq y_t\}} \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] = 0$$



Step 1: bound the average risk

The difference

$$\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(x_t) \neq y_t\}}$$

is a **martingale difference sequence** because

$$\mathbb{E} \left[\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(x_t) \neq y_t\}} \mid (x_1, y_1), \dots, (x_{t-1}, y_{t-1}) \right] = 0$$

The associated martingale is

$$\begin{aligned} & \sum_{t=1}^n \left(\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(x_t) \neq y_t\}} \right) \\ & \iff \underbrace{\frac{1}{n} \sum_{t=1}^n \text{risk}(H_{t-1})}_{\text{average risk}} - \underbrace{\frac{1}{n} \sum_{t=1}^n \mathbb{I}_{\{H_{t-1}(x_t) \neq y_t\}}}_{\text{fraction of mistakes}} \end{aligned}$$



Bernstein's bound

If Z_1, Z_2, \dots is a **martingale difference sequence** with increments bounded by 1 and

$$V_n = \sum_{t=1}^n \mathbb{E} [Z_t^2 \mid Z_1, \dots, Z_{t-1}]$$

then for all $S, K > 0$

$$\mathbb{P} \left(\sum_{t=1}^n Z_t \geq S, \quad V_n \leq K \right) \leq \exp \left(-\frac{S^2}{2(S/3 + K)} \right)$$



Application of Bernstein's bound

Since $0 \leq \mathbb{I}_{\{H(\mathbf{x}) \neq y\}} \leq 1$,

$$\begin{aligned} \text{VAR} \left[\mathbb{I}_{\{H_{t-1}(\mathbf{x}_t), y_t\}} \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] \\ \leq \mathbb{E} \left[\text{risk}(H_{t-1}) \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] \end{aligned}$$



Application of Bernstein's bound

Since $0 \leq \mathbb{I}_{\{H(\mathbf{x}) \neq y\}} \leq 1$,

$$\begin{aligned} \text{VAR} \left[\mathbb{I}_{\{H_{t-1}(\mathbf{x}_t), y_t\}} \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] \\ \leq \mathbb{E} \left[\text{risk}(H_{t-1}) \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] \end{aligned}$$

Applying Bernstein's gives

$$\frac{1}{n} \sum_{t=1}^n \text{risk}(H_{t-1}) \leq \frac{M_n}{n} + \frac{c}{n} \left(\ln M_n + \sqrt{M_n \ln M_n} \right) \quad \text{w.h.p.}$$

Where $\frac{M_n}{n} = \frac{1}{n} \sum_{t=1}^n \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq Y_t\}}$ is the **fraction of mistakes**

Step 2: pick a good classifier in the ensemble

- Start from the ensemble H_0, H_1, \dots, H_n
- Do the following:
 - ① test each H_t on $(x_{t+1}, y_{t+1}), \dots, (x_n, y_n)$
 - ② pick $\hat{H} = H_{t^*}$ minimizing a **penalized risk estimate**



Step 2: pick a good classifier in the ensemble

- Start from the ensemble H_0, H_1, \dots, H_n
- Do the following:
 - ① test each H_t on $(x_{t+1}, y_{t+1}), \dots, (x_n, y_n)$
 - ② pick $\hat{H} = H_{t^*}$ minimizing a **penalized risk estimate**

Guaranteed bound

$$\text{risk}(\hat{H}) \leq \frac{M_n}{n} + \frac{c}{n} \left((\ln n)^2 + \sqrt{M_n \ln n} \right) \quad \text{w.h.p.}$$

