

TTIC 31250

An Introduction to the Theory of Machine Learning

Avrim Blum

04/11/22

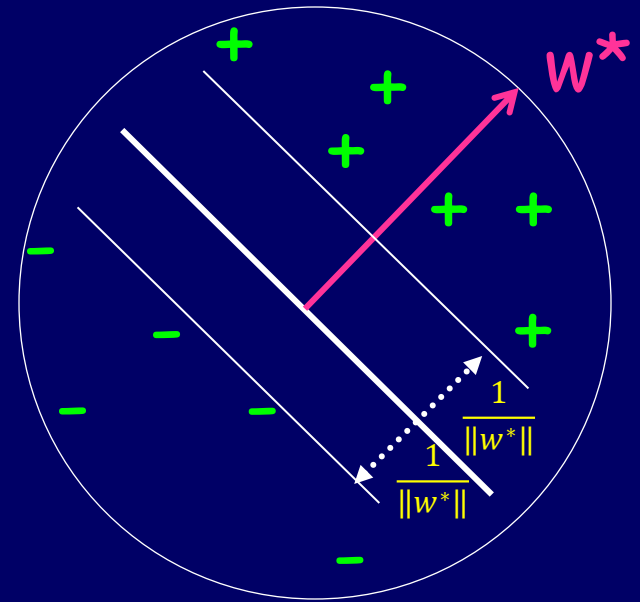
Lecture 4: Support Vector Machines

Perceptron Recap

Perceptron alg makes at most $\|w^*\|^2 R^2$ mistakes if $\exists w^*$ with $w^* \cdot x \geq 1$ on all positives and $w^* \cdot x \leq -1$ on all negatives, and all $\|x\| \leq R$.

Algorithm:

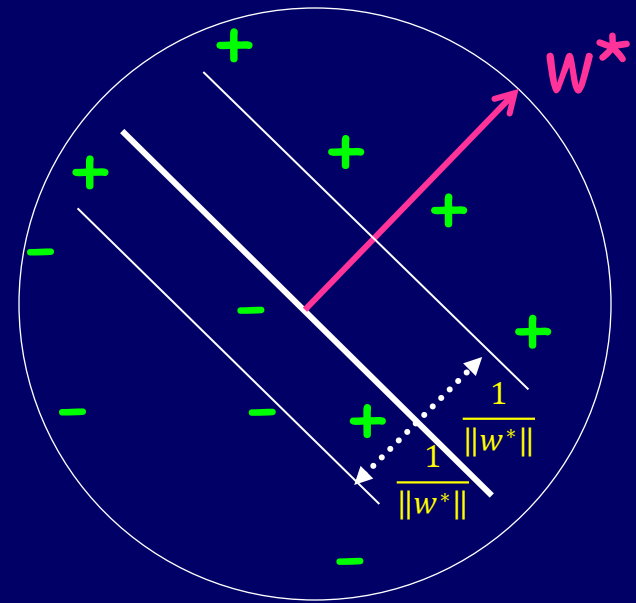
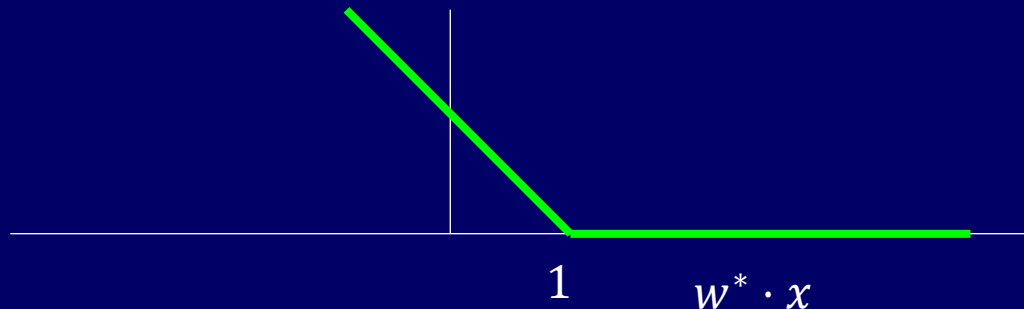
- Initialize $w = \vec{0}$. Predict positive if $w \cdot x > 0$, else predict negative.
- Mistake on positive: $w \leftarrow w + x$.
- Mistake on negative: $w \leftarrow w - x$.



What if w^* isn't perfect?

Theorem: on any sequence of examples S , the Perceptron algo makes at most $\min_{w^*} [\|w^*\|^2 R^2 + 2L_{\text{hinge}}(w^*, S)]$ mistakes.

The **hinge-loss** of w^* on x is the amount by which the desired inequality ($w^* \cdot x \geq 1$ or $w^* \cdot x \leq -1$) is not satisfied.



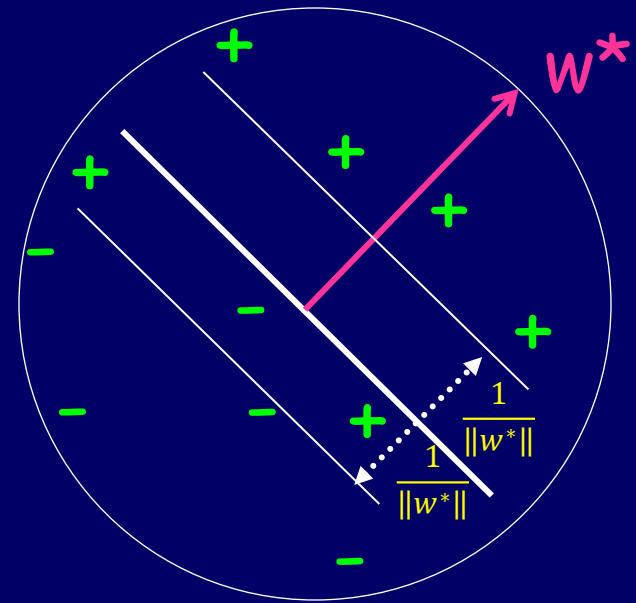
What if w^* isn't perfect?

Theorem: on any sequence of examples S , the Perceptron algo makes at most $\min_{w^*} [\|w^*\|^2 R^2 + 2L_{\text{hinge}}(w^*, S)]$ mistakes.

The **hinge-loss** of w^* on x is the amount by which the desired inequality ($w^* \cdot x \geq 1$ or $w^* \cdot x \leq -1$) is not satisfied.

If x_i is positive then $w^* \cdot x_i \geq 1 - \xi_i$,
if x_i is negative then $w^* \cdot x_i \leq -1 + \xi_i$,
where $\xi_i \geq 0$.

$$\min_{w^*, \xi_1, \xi_2, \dots} [\|w^*\|^2 R^2 + 2 \sum_i \xi_i].$$



What if w^* isn't perfect?

Theorem: on any sequence of examples S , the Perceptron algo makes at most $\min_{w^*} [\|w^*\|^2 R^2 + 2L_{\text{hinge}}(w^*, S)]$ mistakes.

Proof:

$$(w + x_i) \cdot w^* \geq w \cdot w^* + 1 - \xi_i$$

- After M mistakes, $w \cdot w^* \geq M - L_{\text{hinge}}(w^*, S)$.
- Still have: after M mistakes, $\|w\|^2 \leq MR^2$.
- Again use fact that $(w \cdot w^*)^2 \leq \|w\|^2 \|w^*\|^2$.
- Solve: $(M - L_{\text{hinge}})^2 \leq MR^2 \|w^*\|^2$. Do some algebra.

$$M^2 - 2ML_{\text{hinge}} + L_{\text{hinge}}^2 \leq MR^2 \|w^*\|^2$$

$$M \leq R^2 \|w^*\|^2 + 2L_{\text{hinge}} - L_{\text{hinge}}^2/M.$$

Support Vector Machines (SVMs)

In the batch (PAC) setting, we are given S up front. Let's just solve for w^* of largest margin. ("realizable case")

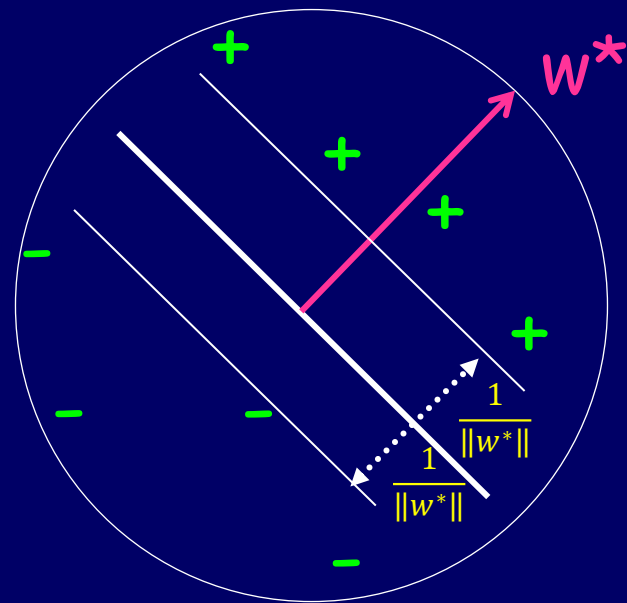
Convex optimization problem:

Minimize: $\|w\|^2$

Subject to: $y_i(w \cdot x_i) \geq 1$ for all $(x_i, y_i) \in S$.

(viewing y_i as ± 1)

But what if there's no perfect separator?



Support Vector Machines (SVMs)

Let's solve for the solution that minimizes a (generalization of) the Perceptron mistake bound.

Given a quantity C as input:

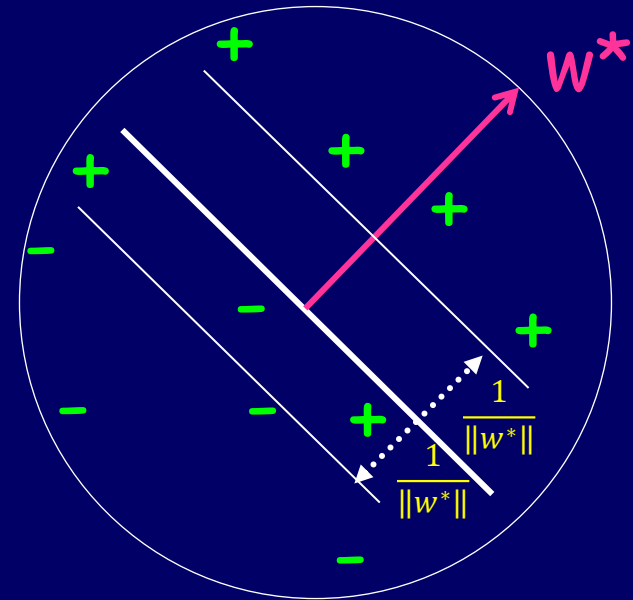
hinge loss. The ξ_i are "slack variables"

$$\text{Minimize: } \|w\|^2 + C \sum_{x_i \in S} \xi_i$$

Subject to: $y_i(w \cdot x_i) \geq 1 - \xi_i$ for all $(x_i, y_i) \in S$.

$$\xi_i \geq 0 \text{ for all } i.$$

This is the SVM algorithm. The quantity C trades off margin and hinge-loss.



Support Vector Machines (SVMs)

Given a quantity C as input:

Minimize: $\|w\|^2 + C \sum_{x_i \in S} \xi_i$

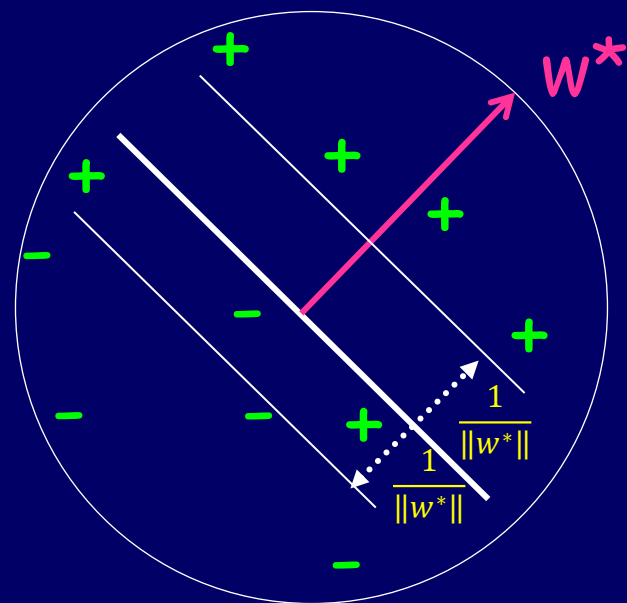
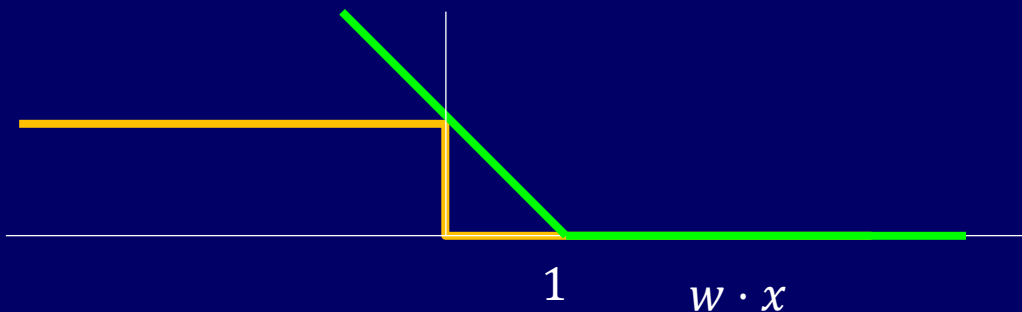
hinge loss. The ξ_i are "slack variables"

Subject to: $y_i(w \cdot x_i) \geq 1 - \xi_i$ for all $(x_i, y_i) \in S$.

$\xi_i \geq 0$ for all i .

Some intuition:

- The total hinge loss is an upper bound on empirical 0/1-loss (# mistakes on S) of the classifier $w \cdot x > 0$.



Support Vector Machines (SVMs)

Given a quantity C as input:

Minimize: $\|w\|^2 + C \sum_{x_i \in S} \xi_i$

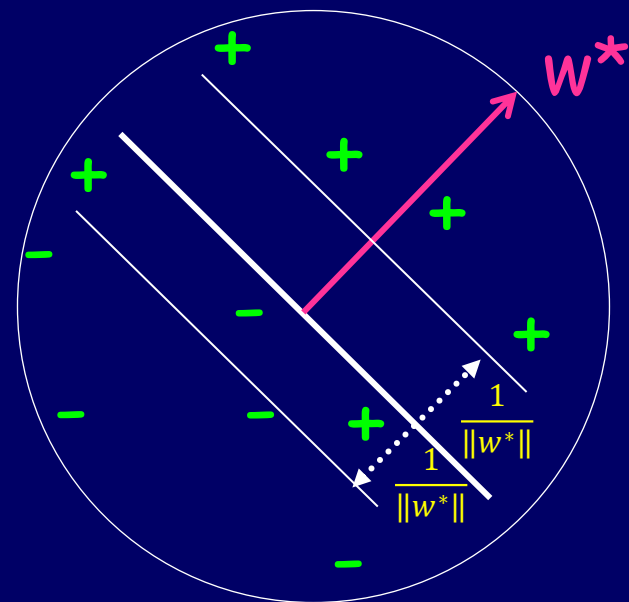
hinge loss. The ξ_i are "slack variables"

Subject to: $y_i(w \cdot x_i) \geq 1 - \xi_i$ for all $(x_i, y_i) \in S$.

$\xi_i \geq 0$ for all i .

Some intuition:

- The total hinge loss is an upper bound on empirical 0/1-loss (# mistakes on S) of the classifier $w \cdot x > 0$.
- The first term $\times R^2$ is roughly (take on faith for now) an upper bound on the amount of overfitting.
- Together, proportional to rough upper-bound on true error.



Support Vector Machines (SVMs)

Given a quantity C as input:

Minimize: $\|w\|^2 + C \sum_{x_i \in S} \xi_i$

hinge loss. The ξ_i are "slack variables"

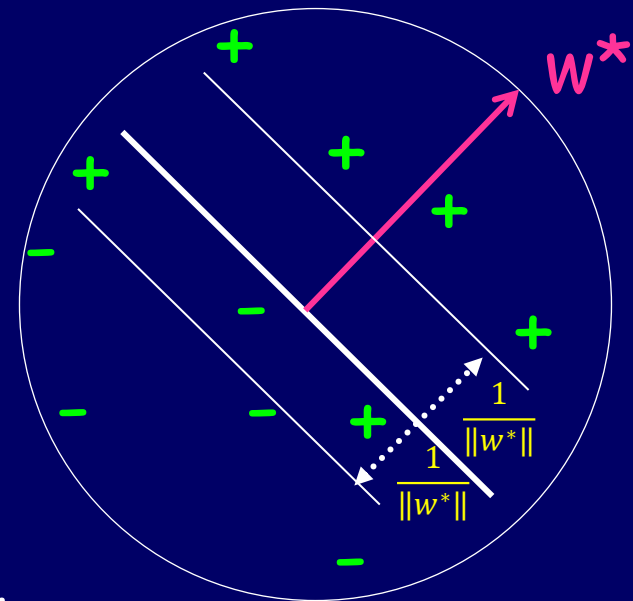
Subject to: $y_i(w \cdot x_i) \geq 1 - \xi_i$ for all $(x_i, y_i) \in S$.

$$\xi_i \geq 0 \text{ for all } i.$$

This is the *primal* form of SVM.

To kernelize it, we will want to move to the *dual* form.

So, first a bit about the Lagrangian dual...



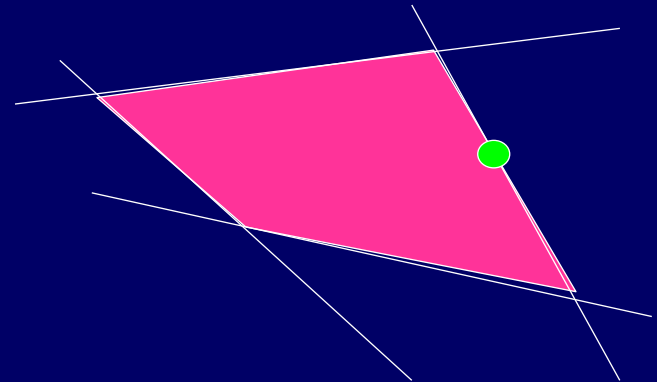
Lagrangian Dual

Consider an optimization problem of the form:

Minimize: convex function in some variables (like w_i, ξ_i)

Subject to: linear constraints on these variables.

Think of as a game between a corporation that wants to minimize its costs (given by the convex function) and a government, that doesn't want the corporation to break any laws (given by the linear constraints).



Lagrangian Dual

Consider an optimization problem of the form:

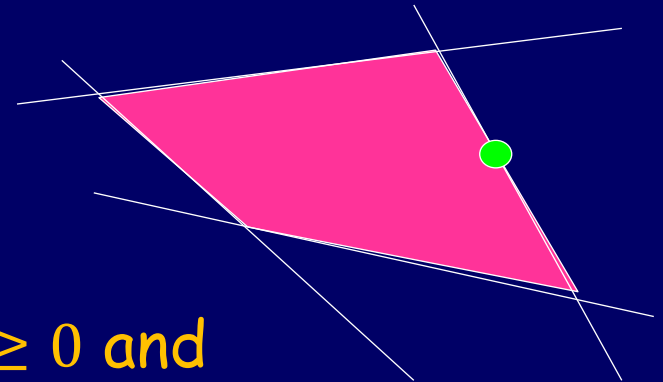
Minimize: convex function in some variables (like w_i, ξ_i)

Subject to: linear constraints on these variables.

For each constraint, the govt can charge a fine that is **linear** in the amount by which it is violated.

E.g., if govt puts fine of \$100 on $\xi_i \geq 0$ and corp uses $\xi_i = -0.5$ then corp pays \$50.

But there's a catch: must be fully linear. If corp uses $\xi_i = +0.5$ then corp collects \$50.



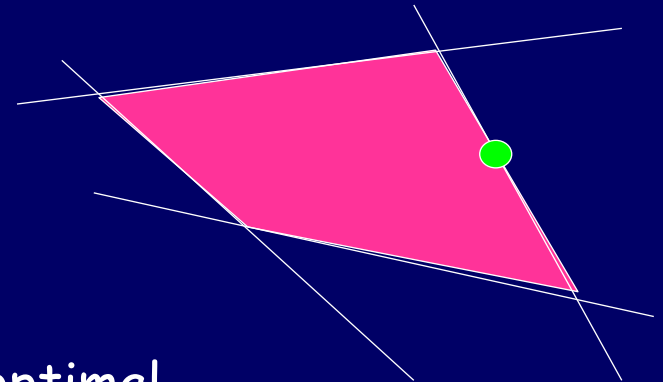
Lagrangian Dual

Consider an optimization problem of the form:

Minimize: convex function in some variables (like w_i, ξ_i)

Subject to: linear constraints on these variables.

The game: for each constraint j , govt gets to choose $\alpha_j \geq 0$ linear penalty. Corp chooses setting of variables. Corp wants to minimize cost. Govt wants to maximize (or equivalently to keep corp honest).



If corp has to go first, clearly should pick optimal feasible point (else govt will assign infinite penalty to any violated constraint).

Claim: If govt goes first, can assign penalties such that corp can do no better. (No "duality gap".) This relies on convexity of the cost function. Govt's optimization problem is called the dual.

Lagrangian Dual

Let \vec{w} denote strategy of corp (primal variables) and let $\vec{\alpha}$ denote strategy of govt (dual variables). The **Lagrangian** is the total cost $L(\vec{w}, \vec{\alpha})$ paid by corp.

Govt's optimization problem is:

$$\max_{\vec{\alpha}} \min_{\vec{w}} L(\vec{w}, \vec{\alpha}) \quad \text{subject to} \quad \alpha_j \geq 0 \quad \forall j$$

Let's see how this plays out for SVMs.

SVM Dual Formulation

Primal: Minimize: $\frac{1}{2} \|w\|^2 + C \sum_i \xi_i$
Subject to: $y_i(w \cdot x_i) \geq 1 - \xi_i$ for all $(x_i, y_i) \in S$.
 $\xi_i \geq 0$ for all i .

Lagrangian: have variables $\alpha_{i1} \geq 0$ for each constraint of 1st kind and $\alpha_{i2} \geq 0$ for each constraint of 2nd kind.

Govt's optimization problem is:

$$\max_{\vec{\alpha}_1, \vec{\alpha}_2} \min_{\vec{w}, \vec{\xi}} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i + \sum_i \alpha_{i1} (1 - \xi_i - y_i(w \cdot x_i)) - \sum_i \alpha_{i2} \xi_i$$

subject to $\alpha_{i1}, \alpha_{i2} \geq 0$ for all i .

SVM Dual Formulation

Now, let's think about a specific ξ_i . Contribution is $\xi_i(C - \alpha_{i1} - \alpha_{i2})$. Govt had better set $\alpha_{i1} + \alpha_{i2} = C$, else corp can make this $-\infty$. So, replace α_{i2} with $C - \alpha_{i1}$, let $\alpha_i = \alpha_{i1}$, and have constraint $0 \leq \alpha_i \leq C$. Simplifies to...

Govt wants to solve for α_i s.t. $0 \leq \alpha_i \leq C$ to maximize

$$\min_w \frac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i(w \cdot x_i))$$

Govt's optimization problem is:

$$\max_{\vec{\alpha}_1, \vec{\alpha}_2} \min_{\vec{w}, \vec{\xi}} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i + \sum_i \alpha_{i1} (1 - \xi_i - y_i(w \cdot x_i)) - \sum_i \alpha_{i2} \xi_i$$

subject to $\alpha_{i1}, \alpha_{i2} \geq 0$ for all i .

SVM Dual Formulation

Now, let's think about a specific ξ_i . Contribution is $\xi_i(C - \alpha_{i1} - \alpha_{i2})$. Govt had better set $\alpha_{i1} + \alpha_{i2} = C$, else corp can make this $-\infty$. So, replace α_{i2} with $C - \alpha_{i1}$, let $\alpha_i = \alpha_{i1}$, and have constraint $0 \leq \alpha_i \leq C$. Simplifies to...

Govt wants to solve for α_i s.t. $0 \leq \alpha_i \leq C$ to maximize

$$\min_w \frac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i (w \cdot x_i))$$

We can solve inner minimization by setting gradient to 0:

$$w - \sum_i \alpha_i y_i x_i = 0.$$

Plug in $w = \sum_i \alpha_i y_i x_i$ above.

SVM Dual Formulation

Dual: solve for α_i s.t. $0 \leq \alpha_i \leq C$ to maximize

$$\begin{aligned} & \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_i \alpha_i - \sum_i \alpha_i y_i \sum_j \alpha_j y_j (x_j \cdot x_i) . \\ & = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) . \end{aligned}$$

Govt wants to solve for α_i s.t. $0 \leq \alpha_i \leq C$ to maximize

$$\min_w \frac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i (w \cdot x_i))$$

We can solve inner minimization by setting gradient to 0:

$$w - \sum_i \alpha_i y_i x_i = 0.$$

Plug in $w = \sum_i \alpha_i y_i x_i$ above.

SVM Dual Formulation

Dual: solve for α_i s.t. $0 \leq \alpha_i \leq C$ to maximize

$$\begin{aligned} & \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_i \alpha_i - \sum_i \alpha_i y_i \sum_j \alpha_j y_j (x_j \cdot x_i) . \\ & = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) . \end{aligned}$$

Notice this is kernelizable. Hence, we can run SVMs with any kernel using the dual and replacing $x_i \cdot x_j$ with $K(x_i, x_j)$.

Intro to Tail Inequalities

Chernoff and Hoeffding bounds

Consider m flips of a coin of bias p . Let N_{heads} be the observed # heads. Let $\epsilon, \alpha \in [0,1]$.

Hoeffding bounds:

- $\Pr[N_{heads}/m > p + \epsilon] \leq e^{-2m\epsilon^2}$
- $\Pr[N_{heads}/m < p - \epsilon] \leq e^{-2m\epsilon^2}$

Chernoff bounds:

- $\Pr[N_{heads}/m > p(1+\alpha)] \leq e^{-mp\alpha^2/3}$
- $\Pr[N_{heads}/m < p(1-\alpha)] \leq e^{-mp\alpha^2/2}$

E.g.,

- $\Pr[N_{heads} > 2(\text{expectation})] \leq e^{-(\text{expectation})/3}$.
- $\Pr[N_{heads} < (\text{expectation})/2] \leq e^{-(\text{expectation})/8}$.

Typical use of bounds

Thm: If $|S| \geq \frac{1}{2\epsilon^2} \left[\ln(2|H|) + \ln\left(\frac{1}{\delta}\right) \right]$, then with prob $\geq 1 - \delta$, all $h \in H$ have $|\text{err}_D(h) - \text{err}_S(h)| < \epsilon$.

- Proof: Just apply Hoeffding + union bound.
 - Chance of failure at most $2|H|e^{-2|S|\epsilon^2}$.
 - Set to δ . Solve.

Hoeffding bounds:

- $\Pr[N_{heads}/m > p + \epsilon] \leq e^{-2m\epsilon^2}$
- $\Pr[N_{heads}/m < p - \epsilon] \leq e^{-2m\epsilon^2}$

Typical use of bounds

Thm: If $|S| \geq \frac{1}{2\epsilon^2} \left[\ln(2|H|) + \ln\left(\frac{1}{\delta}\right) \right]$, then with prob $\geq 1 - \delta$, all $h \in H$ have $|\text{err}_D(h) - \text{err}_S(h)| < \epsilon$.

- Proof: Just apply Hoeffding + union bound.
 - Chance of failure at most $2|H|e^{-2|S|\epsilon^2}$.
 - Set to δ . Solve.
- So, whp, best on sample is ϵ -best over D .
 - Note: this is worse than previous bound ($1/\epsilon$ has become $1/\epsilon^2$), because we are asking for something stronger.
 - Can also get bounds "between" these two.

Typical use of bounds

Thm: If $|S| \geq \frac{6}{\epsilon} \left[\ln |H| + \ln \frac{1}{\delta} \right]$, then with prob $\geq 1 - \delta$,
all $h \in H$ with $\text{err}_D(h) > 2\epsilon$ have $\text{err}_S(h) > \epsilon$, and
all $h \in H$ with $\text{err}_D(h) < \epsilon/2$ have $\text{err}_S(h) < \epsilon$.

Proof: apply Chernoff...

Chernoff bounds:

- $\Pr[N_{heads} / m > p(1+\alpha)] \leq e^{-mp\alpha^2/3}$
- $\Pr[N_{heads} / m < p(1-\alpha)] \leq e^{-mp\alpha^2/2}$

E.g.,

- $\Pr[N_{heads} > 2(\text{expectation})] \leq e^{-(\text{expectation})/3}$.
- $\Pr[N_{heads} < (\text{expectation})/2] \leq e^{-(\text{expectation})/8}$.