



SVMs with Slack & kernels (Not Linearly Separable)

Rishabh Iyer

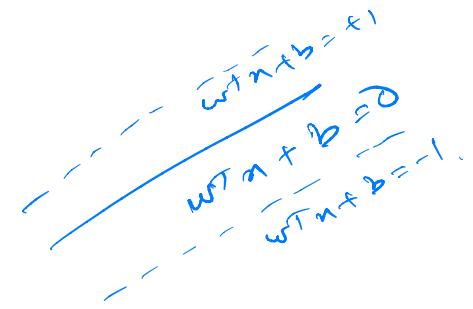
University of Texas at Dallas

Recap of SVM

$$\min_{w,b} \|w\|^2$$

$$\text{s.t. } y_i(w^T n_i + b) \geq 1$$

$$\Rightarrow 1 - y_i(w^T n_i + b) \leq 0 \leftarrow \text{constant}$$



Lagrangian: $L(w, b, \lambda) = \|w\|^2 + \sum_{i=1}^M \lambda_i (1 - y_i(w^T n_i + b))$

$$g(\lambda) = \min_{w,b} L(w,b,\lambda)$$

Dual: $\max_{\lambda: \lambda \geq 0} g(\lambda)$

Dual SVM

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x^{(i)T} x^{(j)} + \sum_i \lambda_i$$

such that

$$\lambda \geq 0 \quad \text{--- } ①$$

$$\sum_i \lambda_i y_i = 0 \quad \text{--- } ②$$

- The dual formulation only depends on inner products between the data points
 - Same thing is true if we use feature vectors instead

Primal

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

$$\text{subject to } y_i(w^T x^{(i)} - b) \geq 1$$

Equivalent:

Dual.

max



$$-\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \gamma_i \gamma_j y_i y_j x^{(i)^T} x^{(j)}$$

s.t.

$$\gamma \geq 0$$

$$\sum_i \gamma_i y_i = 0$$



constraints

$$w = \sum_{i=1}^m \gamma_i y_i x^{(i)}$$

$b \leftarrow \text{complementary slackness}$

Dual SVM

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \Phi(x^{(i)})^T \Phi(x^{(j)}) + \sum_i \lambda_i$$

such that

$$\sum_i \lambda_i y_i = 0$$

Primal:

$$\begin{aligned} \text{min}_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i(w^T \phi(x_i) + b) \geq 1 \end{aligned}$$

- The dual formulation only depends on inner products between the data points
 - Same thing is true if we use feature vectors instead

$$\phi(x) = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$\phi(x) = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

The Kernel Trick

- For some feature vectors, we can compute the inner products quickly, even if the feature vectors are very large
- This is best illustrated by example

- Let $\phi(x_1, x_2) = \begin{bmatrix} x_1 x_2 \\ x_2 x_1 \\ x_1^2 \\ x_2^2 \end{bmatrix}$
- $$\begin{aligned}\phi(x_1, x_2)^T \phi(z_1, z_2) &= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= (x^T z)^2\end{aligned}$$

$$x = [x_1, x_2, \dots, x_n] \in \mathbb{R}^n$$

$$\phi(x) = [\underbrace{x_1, x_2, \dots, x_n}_n, \underbrace{x_1 x_2, \dots}_{{n \choose 2}}, \underbrace{x_1^2, \dots, x_n^2}_n]$$

$$x^T x \approx n^2$$

$$\begin{aligned} & \phi(x_i)^T \phi(x_j) \\ & O(n^2) \rightarrow O(1) \\ & \phi(x_i)^T \phi(x_j) = x_i^T x_j + (x_i^T x_j)^2 \end{aligned}$$

The Kernel Trick

- The same idea can be applied for the feature vector ϕ of all polynomials of degree (exactly) d
 - $\phi(x)^T \phi(z) = (x^T z)^d$
- More generally, a **kernel** is a function $k(x, z) = \phi(x)^T \phi(z)$ for some feature map ϕ
- Rewrite the dual objective

$$\max_{\lambda \geq 0, \sum_i \lambda_i y_i = 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j k(x^{(i)}, x^{(j)}) + \sum_i \lambda_i$$

↑
 Kernel SVM-

Examples of Kernels

- Polynomial kernel of degree exactly d
 - $k(x, z) = \underbrace{(x^T z)}_d \leftarrow$
- General polynomial kernel of degree d for some c
 - $k(x, z) = (x^T z + c)^d$
- Gaussian kernel for some σ
 - $k(x, z) = \exp\left(\frac{-\|x-z\|^2}{2\sigma^2}\right)$
 - The corresponding ϕ is infinite dimensional!
- So many more...

$$\phi(x) = \begin{bmatrix} x \\ x^d \\ \vdots \\ x^{d-1} \end{bmatrix}$$

Gaussian Kernels

- Consider the Gaussian kernel

$$\begin{aligned}\exp\left(\frac{-\|x - z\|^2}{2\sigma^2}\right) &= \exp\left(\frac{-(x - z)^T(x - z)}{2\sigma^2}\right) \\ &= \exp\left(\frac{-\|x\|^2 + 2x^Tz - \|z\|^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right) \exp\left(-\frac{\|z\|^2}{2\sigma^2}\right) \underbrace{\exp\left(\frac{x^Tz}{\sigma^2}\right)}_{\text{red bracket}}\end{aligned}$$

- Use the Taylor expansion for $\exp()$

$$\exp\left(\frac{x^Tz}{\sigma^2}\right) = \sum_{n=0}^{\infty} \frac{(x^Tz)^n}{\sigma^{2n} n!}$$

Kernels



- Bigger feature space increases the possibility of overfitting
 - Large margin solutions may still generalize reasonably well
- Alternative: add “penalties” to the objective to disincentivize complicated solutions

— Regularization ↗
— slack.

Kernel Properties

① if k_1 & k_2 are kernels, is $K = k_1 + k_2$ also a kernel?

$$K(x, y) = \phi(x)^T \phi(y)$$

$$k_1(x, y) = \phi_1(x)^T \phi_1(y)$$

$$k_2(x, y) = \phi_2(x)^T \phi_2(y)$$

$$\phi_3(x) = [\phi_1(x) \ \phi_2(x)]$$

$$\phi_3(y) = [\phi_1(y) \ \phi_2(y)]$$

② if k is a kernel, is $K' = \alpha k$, for $\alpha \geq 0$ also a kernel?

$$\phi_3(x)^T \phi_3(y) = [\phi_1(x)^T \ \phi_2(x)^T] \begin{bmatrix} \phi_1(y) \\ \phi_2(y) \end{bmatrix} = \phi_1(x)^T \phi_1(y) + \phi_2(x)^T \phi_2(y) = k_1 + k_2$$

Kernel SVM

Dual: $\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j k(x^{(i)}, x^{(j)}) + \sum_i \lambda_i$

s.t. $\sum_i \lambda_i y_i = 0, \lambda \geq 0$

Prediction problem?

$$y_t = \text{Sign} \left(\sum_i \lambda_i y_i k(x^{(i)}, x_t) + b \right)$$
$$b = y_i - \sum_j \lambda_j y_j k(x^{(i)}, x^{(j)}) \quad \text{for } i: \lambda_i > 0$$

w, b

$$y_t = \text{sign}(w^T \phi(x_t) + b)$$

$$w = \sum_{i=1}^m \gamma_i y_i \phi(x^{(i)})$$

$$y_t = \text{sign}\left(\sum_{i=1}^m \gamma_i y_i \underbrace{\phi(x^{(i)}) \phi(x_t)}_{K(x^{(i)}, x_t)} + b\right)$$

$$= \text{Sign}\left(\sum_{i=1}^m \gamma_i y_i K(x^{(i)}, x_t) + b\right)$$

Comp. slackness: $b = y_i - w \cdot x_i$, for any $i: \gamma_i \geq 0$

$$= y_i - \sum_{j=1}^m \gamma_j y_j \underbrace{\phi(x^{(j)}) \phi(x^{(i)})}_{K(x^{(j)}, x^{(i)})}$$

SVMs with Slack (Remove Linear Separability)



- Allow misclassification
 - Penalize misclassification linearly (just like in the perceptron algorithm)
 - Again, easier to work with than counting misclassifications
 - Objective stays convex
 - Will let us handle data that isn't linearly separable!
 - Idea: Take the constraints into the main objective
 - The objective function then becomes exactly like what we have seen in Perceptron/Linear Regression

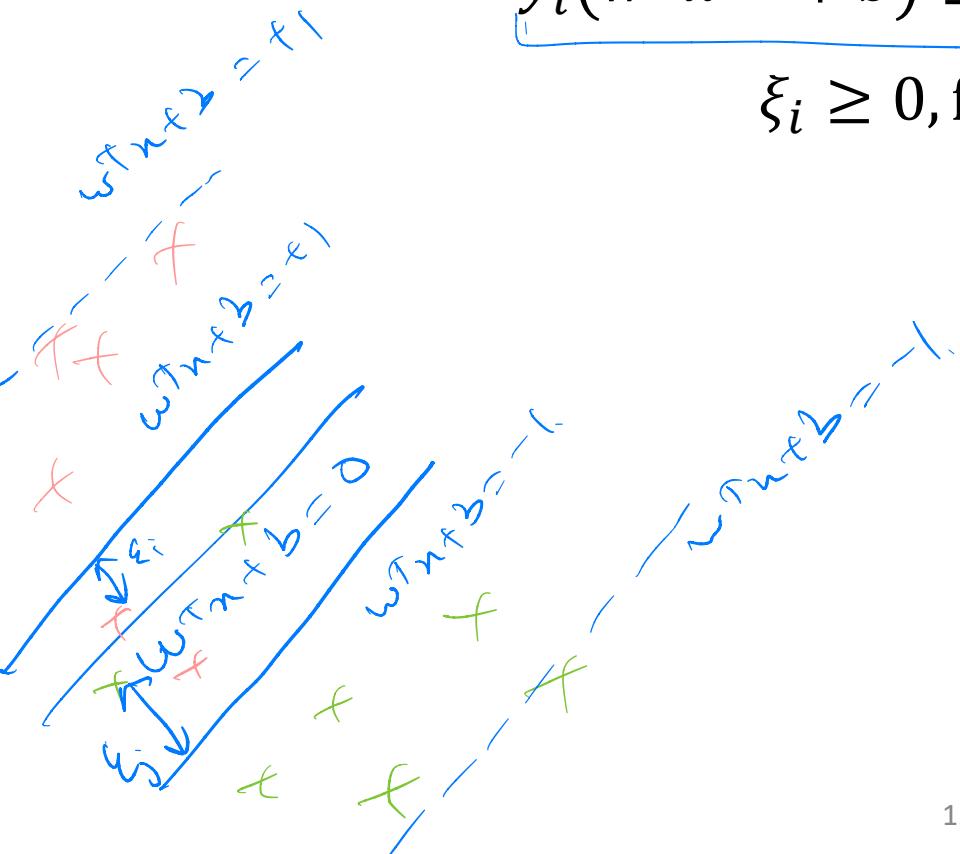
SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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



SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

- How does this objective change with c ?

SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

$c \rightarrow \infty$
 $c = 10^{10}$
 $c = -\infty$
 $\xi_i = 10$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i \text{ for all } i$$

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

- How does this objective change with c ?
 - As $c \rightarrow \infty$, requires a perfect classifier
 - As $c \rightarrow 0$, allows arbitrary classifiers (i.e., ignores the data)

SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

Hyper-parameters

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

- How should we pick c ?

SVMs with Slack



$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

Train / Test

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

Train → 100%
Test → 30%

- How should we pick c ? [hyper-parameters OPT]
- Divide the data into three pieces training, testing, and validation
- Use the validation set to tune the value of the hyperparameter c

Evaluation Methodology

- General learning strategy
 - Build a classifier using the training data
 - Select hyperparameters using validation data →
 - Evaluate the chosen model with the selected hyperparameters on the test data

How can we tell if we overfit the training data?

ML in Practice



- Gather Data + Labels
 - Select feature vectors
 - Randomly split into three groups
 - Training set
 - Validation set
 - Test set
 - Experimentation cycle
 - Select a “good” hypothesis from the hypothesis space
 - Tune hyper-parameters using validation set
 - Compute accuracy on test set (fraction of correctly classified instances)
-
- train-
↓
 $w^{ntn} \rightarrow w_c^*, b_c^*$
↑
tune c
Train → w_c^*, b_c^*

SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i (\xi_i) \leftarrow \# \text{ misClass}^{\curvearrowleft}$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

- What is the optimal value of ξ for fixed w and b ?

$$\xi_i = 1 - y_i (w^T x^{(i)} + b) , \text{ if } y_i - \leq 1 \\ \text{else}$$

$$\xi_i = \max(0, 1 - y_i (w^T x^{(i)} + b))$$

SVMs with Slack

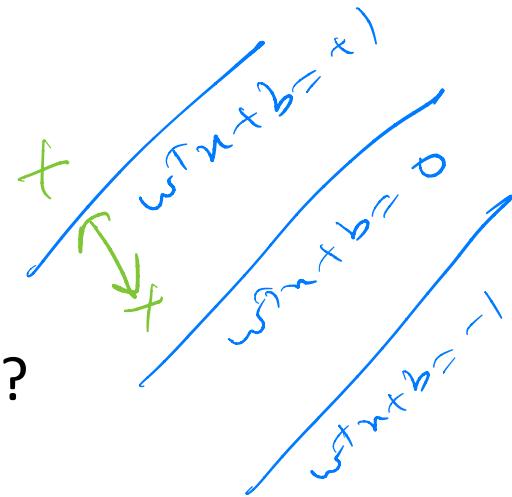
$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

- What is the optimal value of ξ for fixed w and b ?
 - If $y_i(w^T x^{(i)} + b) \geq 1$, then $\xi_i = 0$
 - If $y_i(w^T x^{(i)} + b) < 1$, then $\xi_i = 1 - y_i(w^T x^{(i)} + b)$



SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

- We can formulate this slightly differently

- $\xi_i = \max\{0, 1 - y_i(w^T x^{(i)} + b)\}$

- Does this look familiar?
- Hinge loss provides an upper bound on Hamming loss

Hinge Loss Formulation

- Obtain a new objective by substituting in for ξ

$$\min_{w,b} \frac{1}{2} \|w\|^2 + c \sum_i \max\{0, 1 - y_i(w^T x^{(i)} + b)\}$$

Can minimize with gradient descent!

Hinge Loss Formulation

- Obtain a new objective by substituting in for ξ

$$\min_{w,b} \frac{1}{2} \|w\|^2 + c \sum_i \max\{0, 1 - y_i(w^T x^{(i)} + b)\}$$

Reg
Penalty to prevent overfitting

Loss
Hinge loss

$\min_{w,b} \sum_{i=1}^M L(f(w,b, x^{(i)}), y^{(i)}) + \frac{\lambda}{2} \|w\|^2$

$L(f(w,b, x^{(i)}), y^{(i)}) = \max\{0, 1 - f(w,b, x^{(i)})\}$

$\lambda = \frac{1}{c}$

$$\text{Hinge Loss} = C \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) + \frac{1}{2} \|\mathbf{w}\|^2$$

C is Large; focus on sep/loss
(\rightarrow train error)

\rightarrow not gen.

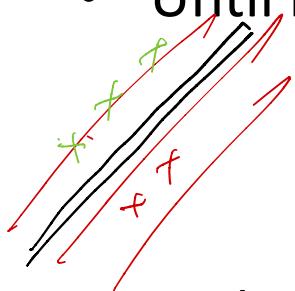
C is small/reasonable gen-

$C \rightarrow 0$ / Arb classifier
Do not care about
loss | Data

REGULARIZATION!!!!



- Until now, we have seen the following optimization problems:



$$\min_{w,b} \sum_i L(f(x^{(i)}, w, b), y_i)$$



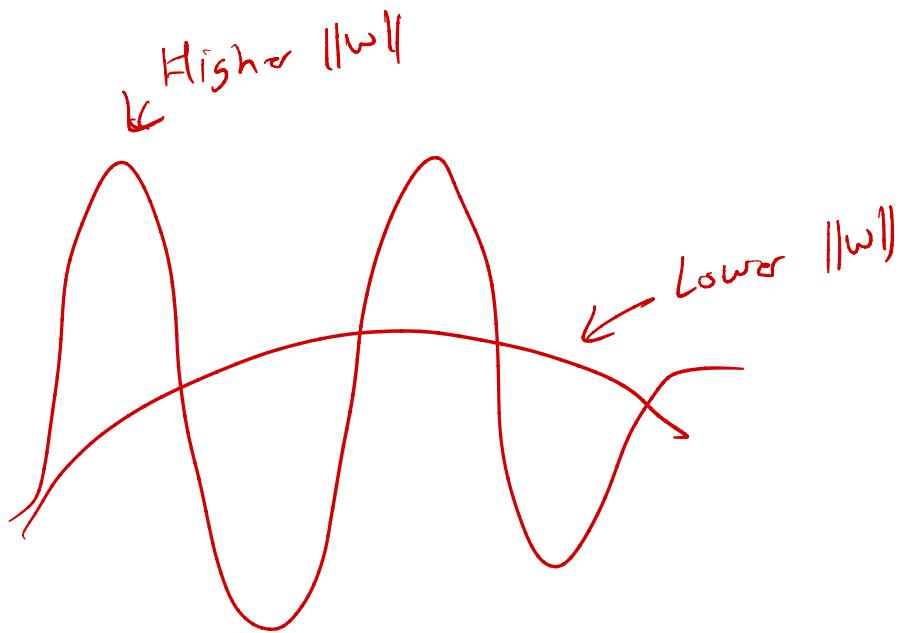
- In the case of Linear regression, L was the squared loss
- In Perceptron, L was Perceptron Loss
- The regularized version of this is:

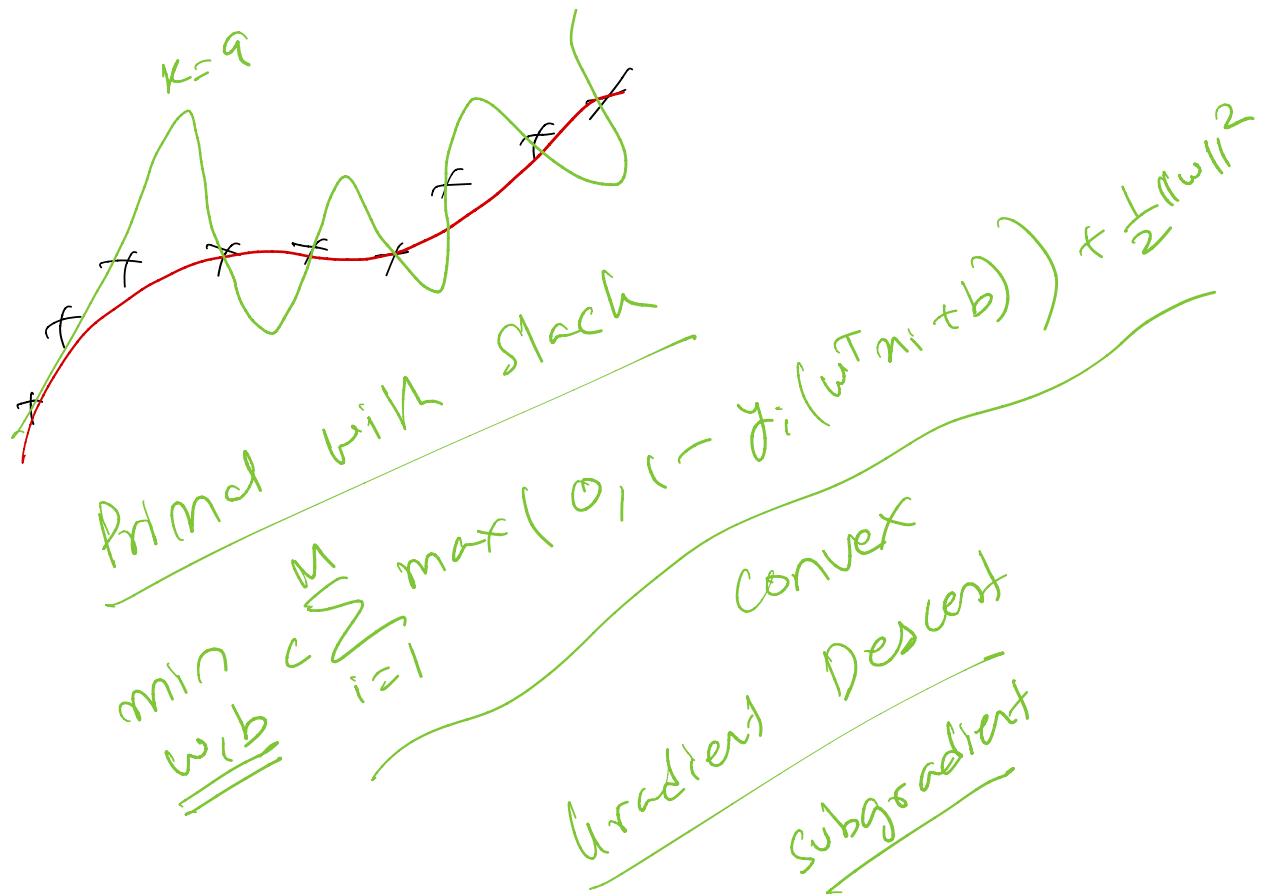
$$\min_{w,b} \frac{1}{2} \cancel{\|w\|^2} + c \sum_i L(f(x^{(i)}, w, b), y_i)$$

- c is a hyper-parameter (again, to be tuned on validation set)

$$\min_{w,b} \sum_{i=1}^m L(f(x^{(i)}, w, b), y^{(i)}) + \frac{c}{2} \|w\|^2$$

$$\min_{\omega, b} \sum_{i=1}^M [(y_i - (\omega^T x_i + b)]^2 + \lambda \|\omega\|^2$$





Perceptron vs Hinge vs Square vs Zero-One Loss



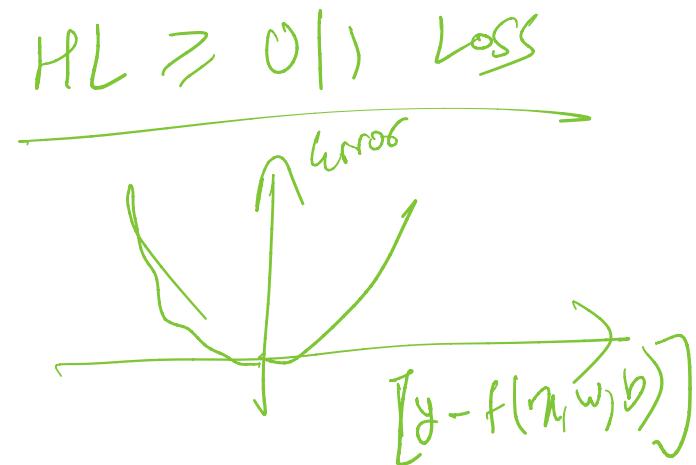
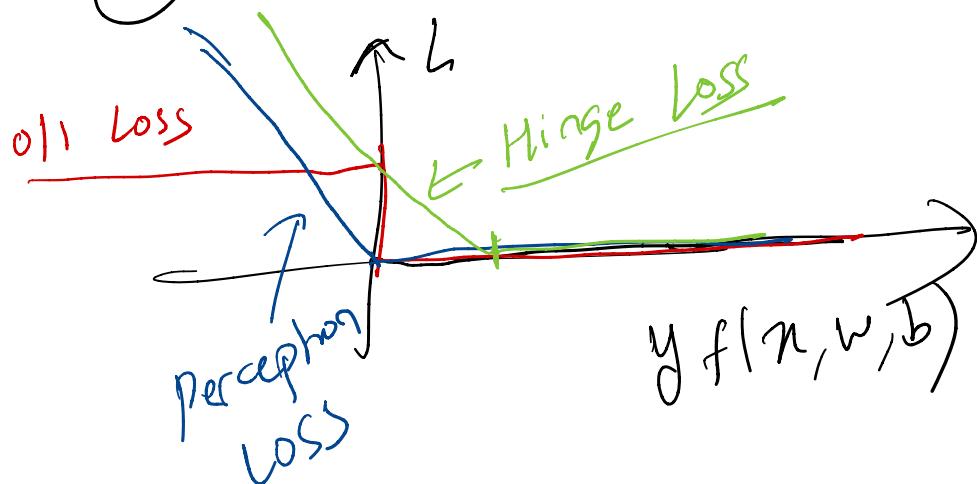
$$\textcircled{1} \quad L(f(x, w, b), y) = \begin{cases} 1 & \text{sign}(f(x, w, b)) \neq y \\ 0 & \text{otherwise} \end{cases}$$

$f(w, b, x) = w^T x + b$

$\xrightarrow{\text{0/1 Loss}} = \frac{1}{2} |\text{sign}(f(x, w, b)) - y|$

$$\textcircled{2} \quad L(f(x, w, b), y) = \max(0, -y f(x, w, b))$$

$$\textcircled{3} \quad L(f(x, w, b), y) = \max(0, 1 - y f(x, w, b))$$



Imbalanced Data



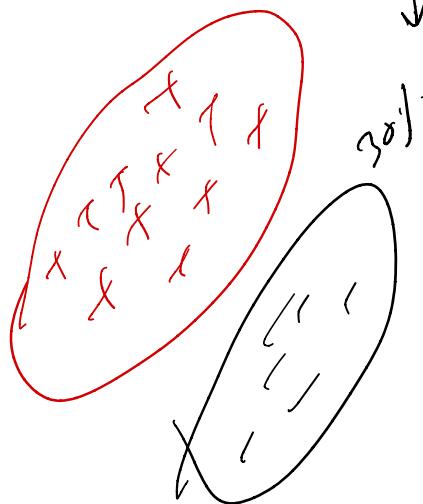
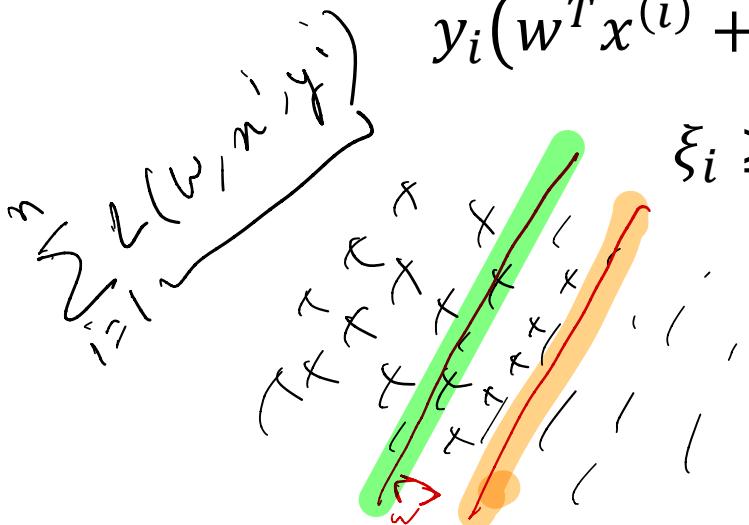
- If the data is imbalanced (i.e., more positive examples than negative examples), may want to evenly distribute the error between the two classes

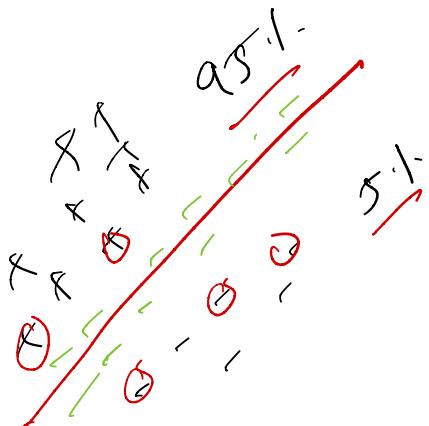
$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + \frac{c}{N_+} \sum_{i:y_i=1} \xi_i + \frac{c}{N_-} \sum_{i:y_i=-1} \xi_i$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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





SVM

$$\min \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1$$

$$x \sum_{i:y_i=1} \ln(f^{(m)}, w, b)$$

$$x \sum_{i:y_i=-1} \ln(f^{(m)}, w, b)$$

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i:y_i \neq 0}$$



Dual of Slack Formulation



$$\min_w f_8(w)$$

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

$$f_i(w) \leq 0$$

$$h_i(w) \geq 0$$

such that

$$y_i(w^T x^{(i)} + b) \geq 1 - \xi_i, \text{ for all } i$$

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

$$1 - \xi_i - y_i(w^T x^{(i)} + b) \leq 0 \quad \leftarrow \text{eq. } i$$

$$-c_i \leq 0 \quad \leftarrow \text{eq. } M_i$$

\neg
ineq.

Dual of Slack Formulation

$$L(w, b, \xi, \lambda, \mu) = \frac{1}{2} w^T w + c \sum_i \xi_i + \sum_i \lambda_i (1 - \xi_i - y_i (w^T x^{(i)} + b)) + \sum_i -\mu_i \xi_i$$

$\sum_i \xi_i [c - \lambda_i - \mu_i] = 0$

Convex in w, b, ξ , so take derivatives to form the dual

$$\begin{cases} w_k \in \mathbb{R} \\ x_{(i)} \in \mathbb{R}^n \end{cases}$$

$$\frac{\partial L}{\partial w_k} = w_k + \sum_i -\lambda_i y_i x_k^{(i)} = 0 \quad \rightarrow \quad w = \sum_i \beta_i y_i x^{(i)}$$

$$\frac{\partial L}{\partial b} = \sum_i -\lambda_i y_i = 0 \quad \Rightarrow \quad \sum_i \beta_i y_i = 0$$

$$\frac{\partial L}{\partial \xi_k} = c - \lambda_k - \mu_k = 0 \quad \Rightarrow \quad \begin{aligned} c &= \beta_k + \mu_k \\ \beta_k &\leq c \end{aligned}$$

$$\begin{matrix} \beta_k \geq 0 \\ \mu_k \geq 0 \end{matrix}$$

Dual of Slack Formulation



$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x^{(i)T} x^{(j)} + \sum_i \lambda_i$$



$$(\omega^T n_i + b)$$

such that

such that

$$\frac{1}{2} \|w\|^2 + C \sum_i \max(0, 1 - y_i w^\top b_i)$$

$$\sum_i \lambda_i y_i = 0$$

$$c \geq \lambda_i \geq 0, \text{ for all } i$$

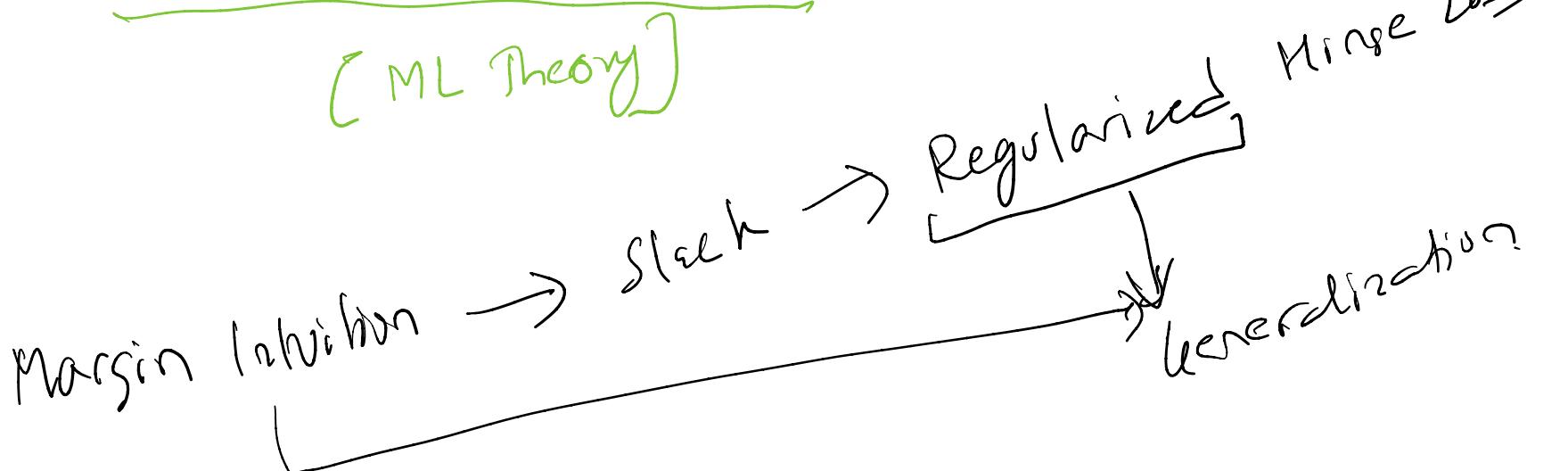
Regularization
in Dual.

$\gamma_i \subseteq c_{\gamma^+ i}$

Generalization



- We argued, intuitively, that SVMs generalize better than the perceptron algorithm → Margin Generalization through Regularization
 - How can we make this precise?



Roadmap

- Where are we headed?
 - Other simple hypothesis spaces for supervised learning
 - k nearest neighbor
 - Decision trees
 - Learning theory
 - Generalization and PAC bounds
 - VC dimension
 - Bias/variance tradeoff



Naive Bayes
Logistic regression

