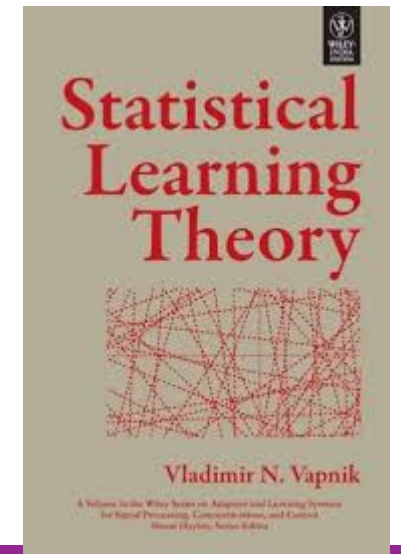


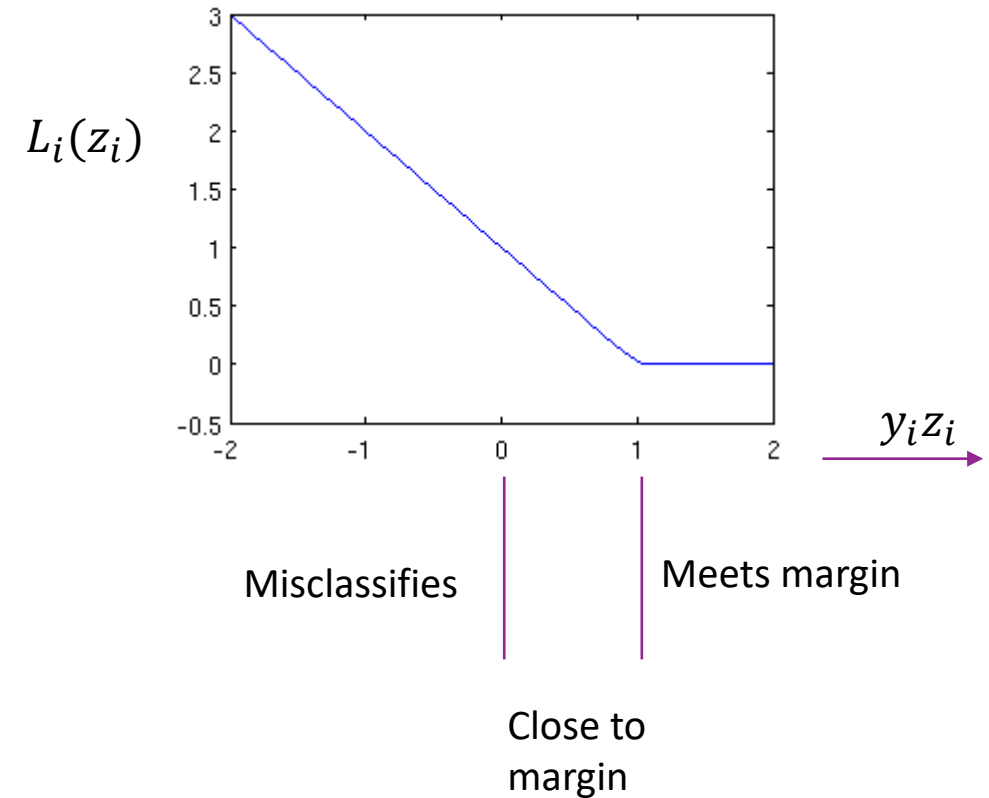
Support Vector Machine

- ❑ Support Vector Machine (SVM)
 - Vladimir Vapnik, 1963
 - But became widely-used with kernel trick, 1993
 - More on this later
- ❑ Got best results on character recognition
- ❑ Key idea: Allow “slack” in the classification
 - Support vector classifier (SVC): Directly use raw features. Good when the original feature space is roughly linearly separable
 - Support vector machine (SVM): Map the raw features to some other domain through a kernel function



Hinge Loss

- Fix $\gamma = 1$
- Want ideally: $y_i(\mathbf{w}^T \mathbf{x} + b) \geq 1$ for all samples i
 - Equivalently, $y_i z_i \geq 1$, $z_i = b + \mathbf{w}^T \mathbf{x}$
- But perfect separation may not be possible
- Define **hinge loss** or **soft margin**:
 - $L_i(\mathbf{w}, b) = \max(0, 1 - y_i z_i)$
- Starts to increase as sample is misclassified:
 - $y_i z_i \geq 1 \Rightarrow$ Sample meets margin target, $L_i(\mathbf{w}) = 0$
 - $y_i z_i \in [0, 1) \Rightarrow$ Sample margin too small, small loss
 - $y_i z_i \leq 0 \Rightarrow$ Sample misclassified, large loss



SVM Optimization

□ Given data (\mathbf{x}_i, y_i)

□ Optimization $\min_{w,b} J(\mathbf{w}, b)$

$$J(\mathbf{w}, b) = C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) + \frac{1}{2} \|\mathbf{w}\|^2$$

C controls final margin

Hinge loss term
Attempts to reduce
Misclassifications

margin = $1/\|\mathbf{w}\|$

□ Constant $C > 0$ will be discussed below

□ Note: ISL book uses different naming conventions.

- We have followed convention in sklearn

Alternate Form of SVM Optimization

□ Equivalent optimization:

$$\min J_1(\mathbf{w}, b, \boldsymbol{\epsilon}), \quad J_1(\mathbf{w}, b, \boldsymbol{\epsilon}) = C \sum_{i=1}^N \epsilon_i + \frac{1}{2} \|\mathbf{w}\|^2$$

□ Subject to constraints:

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i \text{ for all } i = 1, \dots, N$$

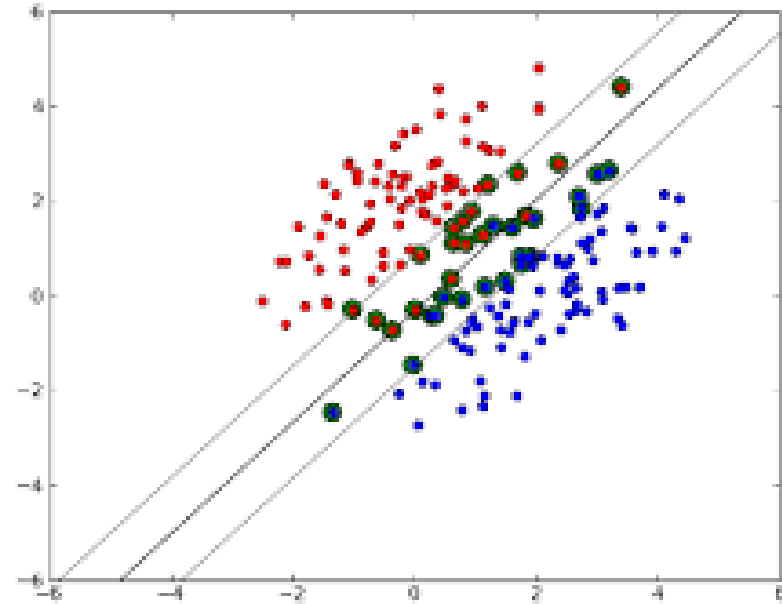
- ϵ_i = amount sample i misses margin target

□ Sometimes write as $J_1(\mathbf{w}, b, \boldsymbol{\epsilon}) = C \|\boldsymbol{\epsilon}\|_1 + \frac{1}{2} \|\mathbf{w}\|^2$

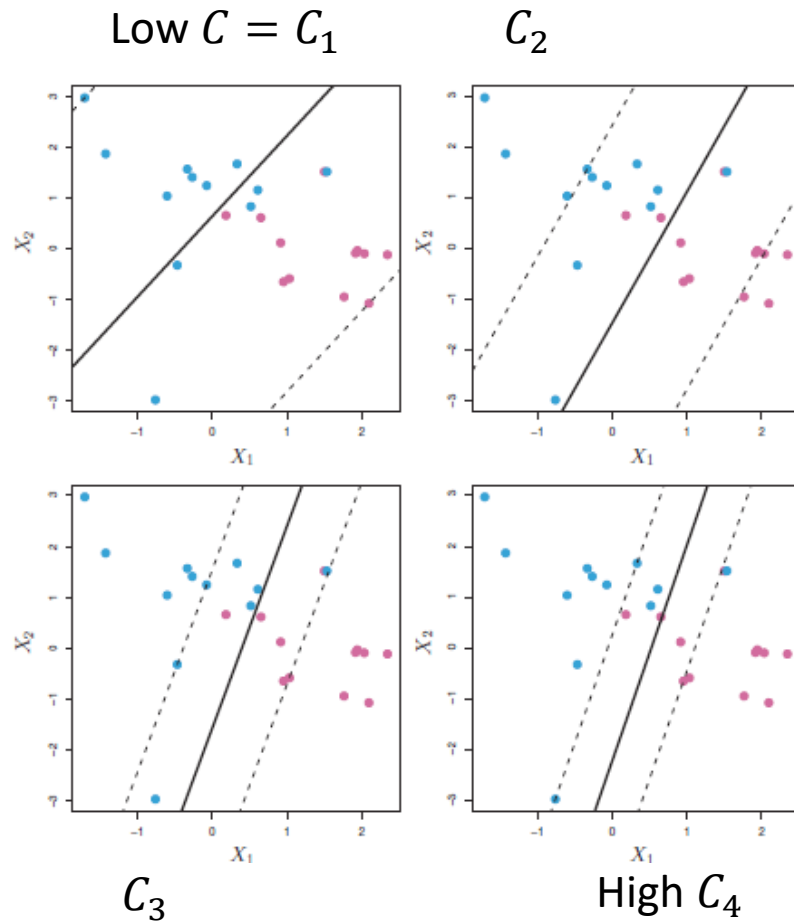
- $\|\boldsymbol{\epsilon}\|_1 = \sum_{i=1}^N \epsilon_i$ called the “one-norm”
- Generally one-norm would have absolute sign over ϵ_i .
- But in this case, when the constraint is met, $\epsilon_i \geq 0$.

Support Vectors

- ❑ **Support vectors:** Samples that either:
 - Are exactly on margin: $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$
 - Or, on wrong side of margin: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 1$
- ❑ Changing samples that are not SVs
 - Does not change solution
 - Provides robustness



Illustrating Effect of C



□ Fig. 9.7 of ISL

- Note: C has opposite meaning in ISL than python
- Here, we use python meaning

□ Low C :

- Leads to large margin
- But allow many violations of margin.
- Many more SVs
- Reduces variance by using more samples

□ Large C :

- Leads to small margin
- Reduce number of violations, and fewer SVs.
- Highly fit to data. Low bias, higher variance
- More chance to overfit

Relation to Logistic Regression

□ Logistic regression also minimizes a loss function:

$$J(\mathbf{w}, b) = \sum_{i=1}^N L_i(\mathbf{w}, b), \quad L_i(\mathbf{w}, b) = \ln P(y_i | \mathbf{x}_i) = -\ln(1 + e^{-y_i z_i})$$

